

## Classification

6 : Predict Bank Credit Risk : Project Title Predict Bank Credit Risk using South German Credit Data Domain Banking

Problem Statement: Normally, most of the bank's wealth is obtained from providing credit loans so that a marketing bank must be able to reduce the risk of non-performing credit loans. The risk of providing loans can be minimized by studying patterns from existing lending data. One technique that you can use to solve this problem is to use data mining techniques. Data mining makes it possible to find hidden information from large data sets by way of classification. The goal of this project, you have to build a model to predict whether the person, described by the attributes of the dataset, is a good (1) or a bad (0) credit risk

Dataset Link : <https://www.kaggle.com/competitions/south-german-credit-prediction/data> (<https://www.kaggle.com/competitions/south-german-credit-prediction/data>)

### Dataset Fields

The below list consists of a detailed breakdown of all the features in the Dataset. 'Kredit' is our target variable, the one whose value must be predicted. P.S. The feature names are in German to preserve the authenticity of the data. In the dataset, they can be replaced using the DataDescription.csv file and the Pandas library.

Id = Id of individual entries, for evaluation

laufkont = status

1 : no checking account

2 : ... < 0 DM 3 : 0 ≤ ... < 200 DM 4 : ... ≥ 200 DM / salary for at least 1 year

laufzeit = duration

moral = credit\_history

0 : delay in paying off in the past

1 : critical account/other credits elsewhere

2 : no credits taken/all credits paid back duly

3 : existing credits paid back duly till now

4 : all credits at this bank paid back duly

verw = purpose

0 : others

1 : car (new)

2 : car (used)

3 : furniture/equipment

4 : radio/television

5 : domestic appliances

6 : repairs

7 : education

8 : vacation

9 : retraining

10 : business

hoehe = amount

sparkont = savings

1 : unknown/no savings account

```
2 : ... < 100 DM 3 : 100 <= ... < 500 DM 4 : 500
<= ... < 1000 DM 5 : ... >= 1000 DM
```

```
beszeit = employment_duration
```

```
1 : unemployed
```

```
2 : < 1 yr 3 : 1 <= ... < 4 yrs 4 : 4 <= ... < 7
yrs 5 : >= 7 yrs
```

```
rate = installment_rate
```

```
1 : >= 35
```

```
2 : 25 <= ... < 35
```

```
3 : 20 <= ... < 25
```

```
4 : < 20
```

```
famges = personal_status_sex
```

```
1 : male : divorced/separated
```

```
2 : female : non-single or male : single
```

```
3 : male : married/widowed
```

```
4 : female : single
```

```
buerge = other_debtors
```

```
1 : none
```

```
2 : co-applicant
```

```
3 : guarantor
```

```
wohnzeit = present_residence
```

```
1 : < 1 yr 2 : 1 <= ... < 4 yrs 3 : 4 <= ... < 7
yrs 4 : >= 7 yrs
```

```
verm = property
```

```
1 : unknown / no property
```

```
2 : car or other
```

```
3 : building soc. savings agr./life insurance
```

```
4 : real estate
```

alter = age

weitkred = other\_installment\_plans

1 : bank

2 : stores

3 : none

wohn = housing

1 : for free

2 : rent

3 : own

bishkred = number\_credits

1 : 1

2 : 2-3

3 : 4-5

4 : >= 6

beruf = job

1 : unemployed/unskilled - non-resident

2 : unskilled - resident

3 : skilled employee/official

4 : manager/self-empl./highly qualif. employee

pers = people\_liable

1 : 3 or more

2 : 0 to 2

telef = telephone

1 : no

2 : yes (under customer name)

gastarb = foreign worker

In [2]:

```
1 # importing Libraries
2 # importing Pandas Library as pd
3 import pandas as pd
4
5 # importing Numpy Library as np
6 import numpy as np
7
8 # importing matplotlib.pyplot as plt
9 import matplotlib.pyplot as plt
10
11 # imporing seaborn as sns
12 import seaborn as sns
13
```

In [3]:

```
1 # Loading the dataset using pandas module and ass
2 df = pd.read_csv('bankcredit.csv')
3
4 # Printing the dataset
5 df
```

Out[3]:

	Id	laufkont	laufzeit	moral	verw	hoehe	spark
0	0	1	18	4	2	1049	
1	1	1	9	4	0	2799	
2	2	2	12	2	9	841	
3	3	1	12	4	0	2122	
4	5	1	10	4	0	2241	
...	...	...	...	...	...	...	
795	993	1	18	4	0	3966	
796	994	1	12	0	3	6199	
797	997	4	21	4	0	12680	

	Id	laufkont	laufzeit	moral	verw	hoehe	spark
798	998	2	12	2	3	6468	
799	999	1	30	2	2	6350	

In [4]:

```
1 df.columns = ['Id','status', 'duration', 'credit_
```

In [5]:

```
1 # By using tails we are getting last 5 values
2 df.tail(10)
```

Out[5]:

	Id	status	duration	credit_history	purpose	am
790	987	1	12		2	3
791	989	2	24		2	0
792	990	1	18		2	6
793	991	2	24		2	1 1
794	992	1	18		2	1
795	993	1	18		4	0
796	994	1	12		0	3
797	997	4	21		4	0 1
798	998	2	12		2	3
799	999	1	30		2	2

10 rows × 22 columns

In [6]:

```
1 # Information about the Dataset:
2 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 800 entries, 0 to 799
```

```
Data columns (total 22 columns):
```

#	Column	Non-Null C
ount	Dtype	
---	-----	
0	Id	800 non-nu
ll	int64	
1	status	800 non-nu
ll	int64	
2	duration	800 non-nu
ll	int64	
3	credit_history	800 non-nu
ll	int64	
4	purpose	800 non-nu
ll	int64	
5	amount	800 non-nu
ll	int64	
6	savings	800 non-nu
ll	int64	
7	employment_duration	800 non-nu
ll	int64	
8	installment_rate	800 non-nu
ll	int64	
9	personal_status_sex	800 non-nu
ll	int64	
10	other_debtors	800 non-nu
ll	int64	
11	present_residence	800 non-nu
ll	int64	
12	property	800 non-nu
ll	int64	
13	age	800 non-nu
ll	int64	
14	other_installment_plans	800 non-nu
ll	int64	
15	housing	800 non-nu
ll	int64	
16	number_credits	800 non-nu
ll	int64	

```
17 job                        800 non-nu
ll  int64
18 people_liable             800 non-nu
ll  int64
19 telephone                 800 non-nu
ll  int64
20 foreign_worker            800 non-nu
ll  int64
21 credit_risk               800 non-nu
ll  int64
dtypes: int64(22)
memory usage: 127.6 KB
```

In [7]:

```
1 # isnull() method is used to check whether the da
2 # sum() is used to find count of null values
3 print(df.isnull().sum())
4
5 #getting size of the dataset(rows multiplied by c
6 print(df.size)
7
8 # getting rows and columnes in a dataset
9 print(df.shape)
```



```
Id                0
status            0
duration          0
credit_history    0
purpose          0
amount           0
savings          0
employment_duration
```

In [10]:

```
1 # getting all statctical values like mean,median,
2 # for numerical values
3 df.describe()
```

Out[10]:

	Id	status	duration	credit_histo
count	800.000000	800.000000	800.000000	800.000000
mean	478.101250	2.648750	20.496250	2.582500
std	278.883661	1.250931	12.006881	1.099800
min	0.000000	1.000000	4.000000	0.000000
25%	238.750000	1.750000	12.000000	2.000000
50%	472.000000	2.000000	18.000000	2.000000
75%	707.250000	4.000000	24.000000	4.000000
max	999.000000	4.000000	72.000000	4.000000

8 rows × 22 columns

In [11]:

```
1 df.columns
```

Out[11]:

```
Index(['Id', 'status', 'duration', 'credit_history', 'purpose', 'amount',  
      'savings', 'employment_duration', 'installment_rate',  
      'personal_status_sex', 'other_debtors', 'present_residence', 'property',  
      'age', 'other_installment_plans', 'housing', 'number_credits', 'risk'])
```

In [12]:

```
1 df.drop(['Id'], inplace=True, axis=1)
```

In [53]:

```
1 df
```

Out[53]:

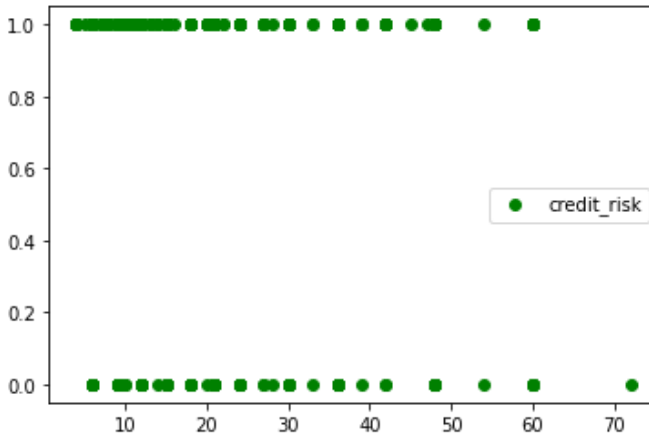
	status	duration	credit_history	purpose	amount
0	1	18	4	2	1049
1	1	9	4	0	2799
2	2	12	2	9	841
3	1	12	4	0	2122
4	1	10	4	0	2241
...	...	...	...	...	...
795	1	18	4	0	3966
796	1	12	0	3	6199
797	4	21	4	0	12680
798	2	12	2	3	6468
799	1	30	2	2	6350

800 rows × 21 columns

## Data Visualization

In [19]:

```
1 #Two column should be numerical
2 plt.scatter(x = 'duration',y = 'credit_risk',data
3
4 #Describe all the elements of a graph
5 plt.legend()
6
7 #show the graph
8 plt.show()
```



In [15]:

```
1 plt.hist(x = df["credit_risk"],color = 'r',bins=5)
2 plt.title("Credit dataset")
3 # plot the grid:
4 plt.grid()
5 # display the plot:
6 plt.show()
```

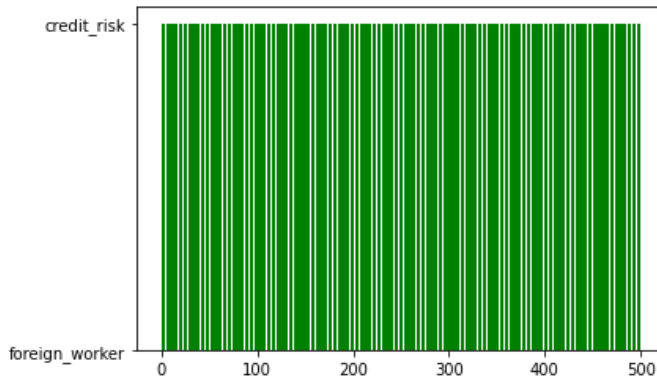


In [16]:

```
1 index = np.arange(500)
2 plt.bar(x = index,height = 'foreign_worker' ,color = 'green')
3 plt.bar(x = index,height = 'credit_risk' ,color = 'red')
```

Out[16]:

<BarContainer object of 500 artists>

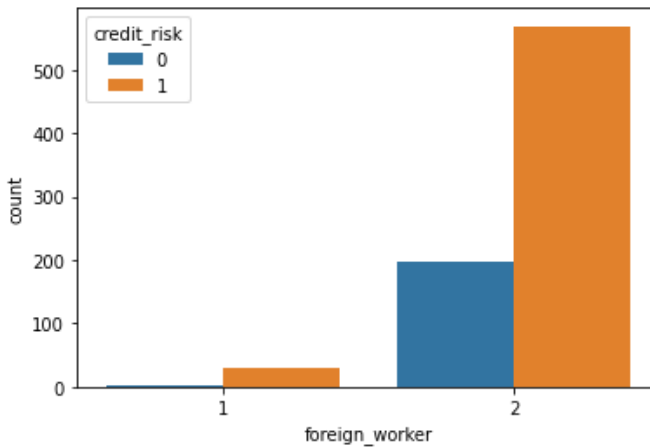


In [17]:

```
1 from warnings import filterwarnings
2 filterwarnings('ignore')
3 sns.countplot(df['foreign_worker'], hue=df['credit_risk'])
```

Out[17]:

<AxesSubplot:xlabel='foreign\_worker', ylabel='count'>

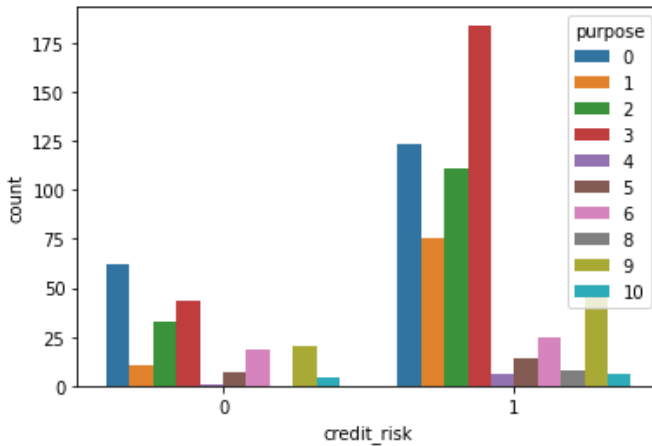


In [18]:

```
1 sns.countplot(df['credit_risk'], hue=df['purpose'])
```

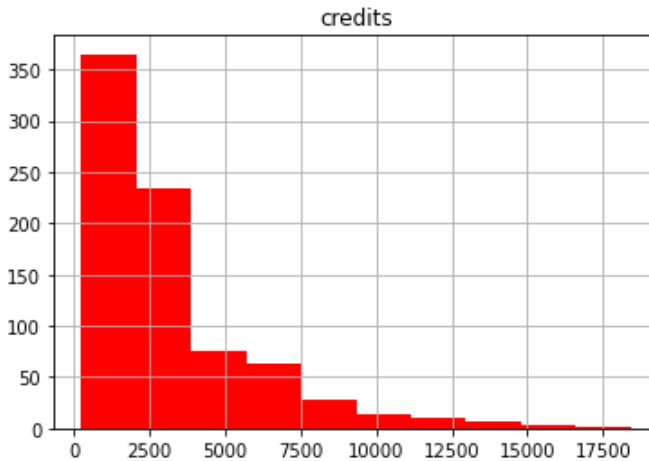
Out[18]:

<AxesSubplot:xlabel='credit\_risk', ylabel='count'>



In [46]:

```
1 #Plot the Histogram with multiple bins and add G
2
3 plt.hist(x = df['amount'],color = 'r')
4 plt.title('credits')
5
6 #Plot the Grid
7 plt.grid()
8 plt.show()
```



In [ ]:

1

In [20]:

```
1 # check the bad loans
2 df[df['credit_risk']==0]
```

Out[20]:

	status	duration	credit_history	purpose	amount
600	1	18	2	0	1216
601	4	18	4	6	1864



	status	duration	credit_history	purpose	amount
602	1	12	2	0	1228
603	2	12	2	0	685
604	3	9	2	3	745
...	...	...	...	...	...
795	1	18	4	0	3966
796	1	12	0	3	6199
797	4	21	4	0	12680
798	2	12	2	3	6468
799	1	30	2	2	6350

In [21]:

```
1 # check the good loans
2 df[df['credit_risk']==1]
```

Out[21]:

	status	duration	credit_history	purpose	amount
0	1	18	4	2	1049
1	1	9	4	0	2799
2	2	12	2	9	841
3	1	12	4	0	2122
4	1	10	4	0	2241
...	...	...	...	...	...
595	4	6	2	2	1766
596	2	24	2	1	2760
597	4	24	4	5	5507
598	3	24	2	2	2892
599	2	36	3	0	2862

600 rows × 21 columns

In [22]:

```
1 df["credit_risk"].value_counts()  
2
```

Out[22]:

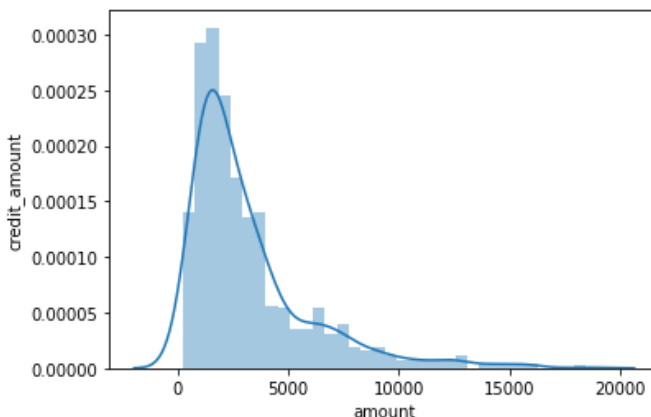
```
1      600  
0      200  
Name: credit_risk, dtype: int64
```

## Data Transformation

In [23]:

```
1 # Distribution of the Dispalcement column  
2 sns.distplot(df['amount'])  
3 plt.ylabel('credit_amount')  
4 print('Skewness : ',df['amount'].skew())  
5 plt.show()
```

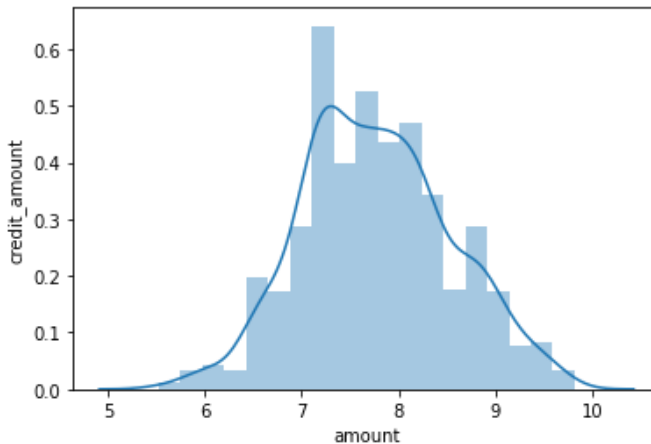
Skewness : 2.0472701844801806



In [24]:

```
1 # Apply natural log transformation for Displaceme
2 log_amount = np.log(df['amount'])
3 print('Skewness after log Transformation : ',log_a
4
5 sns.distplot(log_amount)
6 plt.ylabel('credit_amount')
7 plt.show()
```

Skewness after log Transformation : 0.13  
83727959275712



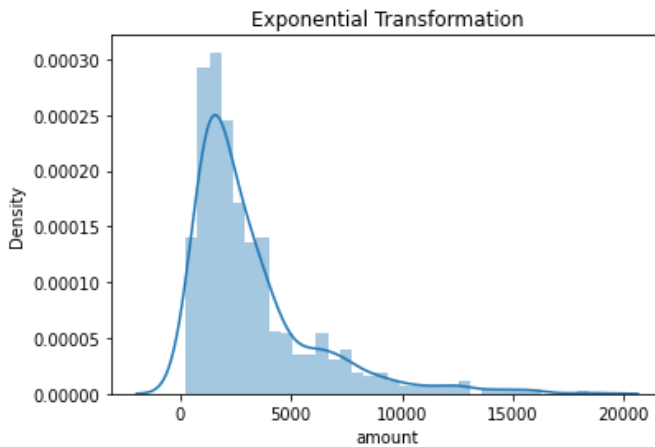
In [25]:

```

1 # Anti - log or Exponential Transformation
2 displacement = np.exp(log_amount)
3 print('Skewness after log Transformation : ',displ
4 # plot the Distribution
5 sns.distplot(displacement)
6 plt.ylabel('Density')
7 plt.title('Exponential Transformation')
8 plt.show()

```

Skewness after log Transformation : 2.04  
727018448018



In [26]:

```

1 df.drop(['age', 'duration'],axis=1,inplace=True)
2 df

```

Out[26]:

	status	credit_history	purpose	amount	savings
0	1	4	2	1049	1
1	1	4	0	2799	1
2	2	2	9	841	2

	status	credit_history	purpose	amount	savings
3	1	4	0	2122	1
4	1	4	0	2241	1
...	...	...	...	...	...
795	1	4	0	3966	1
796	1	0	3	6199	1
797	4	4	0	12680	5
798	2	2	3	6468	5
799	1	2	2	6350	5

In [27]:

```
1 df.head()
```

Out[27]:

	status	credit_history	purpose	amount	savings	ei
0	1	4	2	1049	1	
1	1	4	0	2799	1	
2	2	2	9	841	2	
3	1	4	0	2122	1	
4	1	4	0	2241	1	

In [29]:

```
1 from sklearn.model_selection import train_test_sp
2 from sklearn.metrics import classification_report
```

In [30]:

```
1 xtrain, xtest, ytrain, ytest = train_test_split(d
```

In [31]:

```
1 xtrain.shape, xtest.shape
```

Out[31]:

```
((600, 18), (200, 18))
```

In [33]:

```
1 from sklearn.tree import DecisionTreeClassifier
2 from sklearn.tree import plot_tree
3 dtc = DecisionTreeClassifier()
```

In [34]:

```
1 dtc.fit(xtrain, ytrain)
```

Out[34]:

```
DecisionTreeClassifier()
```

In [35]:

```
1 plot_tree(dtc)
```

Out[35]:

```
[Text(220.88056640625, 211.4, 'X[0] <=
```

In [36]:

```
1 dtc.score(xtrain, ytrain)
```

Out[36]:

1.0

In [37]:

```
1 v_pred = dtc.predict(xtest)
```

In [38]:

```
1 print(classification_report(ytest, v_pred))
```

		precision	recall	f1-score
0	44	0.42	0.52	0.46
1	156	0.86	0.79	0.82
accuracy	200			0.73
macro avg	200	0.64	0.66	0.64
weighted avg	200	0.76	0.73	0.74

In [39]:

```
1 dtc.score(xtest, ytest)
```

Out[39]:

0.735

In [29]:

```
1 DecisionTreeClassifier?
```

In [41]:

```
1 from sklearn.model_selection import GridSearchCV
```

In [42]:

```
1 params = {  
2     'criterion' : ["gini", "entropy", "log_loss"  
3     "splitter" : ["best", "random"],  
4     "max_features" : ["auto", "sqrt", "log2"],  
5     "min_samples_leaf" : [1,2,5,6]  
6 }
```

In [43]:

```
1 grid = GridSearchCV(estimator = dtc, param_grid =
```

In [34]:

```
1 grid.fit(xtrain,ytrain)
```



```
/home/student/my_project/env/lib/python
```

In [44]:

```
1 dtc.score(xtrain, ytrain)
```

Out[44]:

1.0

In [45]:

```
1 dtc.score(xtest, ytest)
```

Out[45]:

0.735

In [46]:

```
1 dtc2 = DecisionTreeClassifier(
2     criterion= 'gini',
3     max_features = 'sqrt',
4     splitter = 'random',
5     min_samples_leaf = 6)
```

In [47]:

```
1 dtc2.fit(xtrain, ytrain)
```

Out[47]:

```
DecisionTreeClassifier(max_features='sqrt', min_samples_leaf=6,
                        splitter='random',
                        ...)
```

In [48]:

```
1 df.describe()
```

Out[48]:

	status	credit_history	purpose	am
count	800.000000	800.000000	800.000000	800.000000

	status	credit_history	purpose	am
mean	2.648750	2.582500	2.785000	3210.290
std	1.250931	1.099866	2.680533	2792.840
min	1.000000	0.000000	0.000000	250.000
25%	1.750000	2.000000	1.000000	1364.000
50%	2.000000	2.000000	2.000000	2264.000
75%	4.000000	4.000000	3.000000	3907.250

In [53]:

```
1 from sklearn.preprocessing import MinMaxScaler
2 min = MinMaxScaler()
```

In [54]:

```
1 for column in df.columns:
2     df[column] = min.fit_transform(df[[column]])
```

In [55]:

```
1 df.describe()
```

Out[55]:

	status	credit_history	purpose	amou
count	800.000000	800.000000	800.000000	800.000000
mean	0.549583	0.645625	0.278500	0.162800
std	0.416977	0.274967	0.268053	0.153600
min	0.000000	0.000000	0.000000	0.000000
25%	0.250000	0.500000	0.100000	0.061200
50%	0.333333	0.500000	0.200000	0.110800
75%	1.000000	1.000000	0.300000	0.201200
max	1.000000	1.000000	1.000000	1.000000

In [60]:

```
1 from sklearn.preprocessing import StandardScaler
2 s = StandardScaler()
```

In [62]:

```
1 selected = pd.DataFrame(s.fit_transform(df.drop("
2 selected["ls"] = df["credit_risk"]
```

In [112]:

```
1 selected["ls"] = le.fit_transform(selected["ls"])
```

In [113]:

```
1 selected
```

Out[113]:

	status	credit_history	purpose	amount	s
0	-1.318842	1.289599	-0.278193	-1.262349	-0.7
1	-1.318842	1.289599	-1.099425	0.316359	-0.7
2	-0.518938	-0.529941	2.185506	-1.428529	-0.0
3	-1.318842	1.289599	-1.099425	-0.140635	-0.7
4	-1.318842	1.289599	-1.099425	-0.052929	-0.7
...	...	...	...	...	...
795	-1.318842	1.289599	-1.099425	0.893372	-0.7
796	-1.318842	-2.349482	0.132424	1.285741	-0.7
797	1.080871	1.289599	-1.099425	1.719655	1.7
798	-0.518938	-0.529941	0.132424	1.350367	1.7
799	-1.318842	-0.529941	-0.278193	1.331902	1.7

800 rows × 19 columns

In [131]:

```
1 from sklearn.linear_model import LogisticRegression
2 from sklearn.tree import DecisionTreeClassifier
3 from sklearn.svm import SVC
4 from sklearn.naive_bayes import GaussianNB
```

In [109]:

```
1 lr = LogisticRegression()
2 dtc = DecisionTreeClassifier()
3 sv = SVC()
4 gn timer = GaussianNB()
```

In [110]:

```
1 g = GridSearchCV(estimator= dtc , param_grid = {})
```

In [118]:

```
1 x = df.drop("credit_risk",axis=1)
2 v = df["credit_risk"]
```

In [119]:

```
1 svc = GaussianNB()
```

In [120]:

```
1 svc.fit(x,v)
```

Out[120]:

GaussianNB()

In [121]:

```
1 svc.score(x,v)
```

Out[121]:

0.76

In [122]:

```
1 x.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 800 entries, 0 to 799
Data columns (total 18 columns):
 #   Column                Non-Null C
ount  Dtype
---  ---
0    status                800 non-nu
ll   int64
1    credit_history        800 non-nu
ll   int64
2    purpose              800 non-nu
ll   int64
3    amount               800 non-nu
ll   int64
4    savings              800 non-nu
ll   int64
5    employment_duration  800 non-nu
ll   int64
6    installment_rate     800 non-nu
ll   int64
7    personal_status_sex  800 non-nu
ll   int64
8    other_debtors        800 non-nu
ll   int64
9    present_residence    800 non-nu
ll   int64
10   property             800 non-nu
ll   int64
11   other_installment_plans 800 non-nu
ll   int64
12   housing              800 non-nu
ll   int64
13   number_credits       800 non-nu
ll   int64
14   job                  800 non-nu
ll   int64
```

```

15  people_liable          800 non-nu
ll   int64
16  telephone            800 non-nu
ll   int64
17  foreign_worker        800 non-nu
ll   int64
dtypes: int64(18)
memory usage: 112.6 KB

```

In [123]:

```
1 svc?
```

In [126]:

```

1 param_grid = [{'n_estimators': [3, 10, 30], 'max_
2               'learning_rate':[0.3,0.5,0.01,0.1]}
3

```

In [132]:

```
1 grid.fit(x,v)
```

Out[132]:

```

GridSearchCV(estimator=DecisionTreeClass
ifier(),
              param_grid={'criterion': ['
gini', 'entropy', 'log_loss'],
                          'max_features':
['auto', 'sqrt', 'log2'],
                          'min_samples_le
af': [1, 2, 5, 6],
                          'splitter': ['b
est', 'random']})

```

In [134]:

```
1 grid = GridSearchCV(estimator=svc, param_grid =
```

In [141]:

```
1 grid.best_params_
```

```

-----
-----
AttributeError
Traceback (most recent call last)
<ipython-input-141-28c2e4d7952c> in <module>
----> 1 grid.best_params_

```

**AttributeError:** 'GridSearchCV' object has no attribute 'best\_params\_'

In [115]:

```

1 from sklearn.ensemble import RandomForestClassifier
2 rfc = RandomForestClassifier()

```

In [68]:

```

1 rfc.fit(xtrain,ytrain)
2
3 rfc.score(xtrain,ytrain)

```

Out[68]:

1.0

In [69]:

```

1 rfc?

```

In [142]:

```

1 pr = {
2     'n_estimators' : [10,100,1000,600],
3     'criterion':["log_loss", "entropy", "gini"]
4 }

```

In [143]:

```

1 g = GridSearchCV(estimator=rfc, param_grid = pr)

```

In [144]:

```
1 g.fit(x,v)
```

Out[144]:

```
GridSearchCV(estimator=RandomForestClassifier(),  
              param_grid={'criterion': ['  
log_loss', 'entropy', 'gini'],  
                          'n_estimators':  
[10, 100, 1000, 600]})
```

In [145]:

```
1 rfc2 = RandomForestClassifier(criterion = 'log_lo
```

In [148]:

```
1 rfc2.fit(x,y)  
2  
3 rfc2.score(x,v)
```

Out[148]:

1.0