

```
In [2]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
import plotly.express as px

base_credir = pd.read_csv(r'C:\Users\Master\Desktop\Jupyter\IA UDEMY\credir_data.csv')
base_credir

Out[2]:
```

	clientid	income	age	loan	default
0	1	66155.925095	59.017015	8106.532131	0
1	2	34415.153966	48.117153	6564.745018	0
2	3	57317.170063	63.108049	8020.953296	0
3	4	42709.534201	45.751972	6103.642260	0
4	5	66952.688845	18.584336	8770.099235	1
...
1995	1996	59221.044874	48.518179	1926.729397	0
1996	1997	69516.127573	23.162104	3503.176156	0
1997	1998	44311.449262	28.017167	5522.786693	1
1998	1999	43756.056605	63.971796	1622.722598	0
1999	2000	69436.579552	56.152617	7378.833599	0

2000 rows × 5 columns

```
In [3]: base_credir.head()
```

```
Out[3]:
```

	clientid	income	age	loan	default
0	1	66155.925095	59.017015	8106.532131	0
1	2	34415.153966	48.117153	6564.745018	0
2	3	57317.170063	63.108049	8020.953296	0
3	4	42709.534201	45.751972	6103.642260	0
4	5	66952.688845	18.584336	8770.099235	1

```
In [4]: base_credir.tail(8)
```

```
Out[4]:
```

1992	1993	30803.806165	23.250084	623.024153	0
1993	1994	54421.410155	26.821928	3273.631823	0
1994	1995	24254.700791	37.751622	2225.284643	0
1995	1996	59221.044874	48.518179	1926.729397	0
1996	1997	69516.127573	23.162104	3503.176156	0
1997	1998	44311.449262	28.017167	5522.786693	1
1998	1999	43756.056605	63.971796	1622.722598	0
1999	2000	69436.579552	56.152617	7378.833599	0

```
In [5]: base_credir.describe()
```

```
Out[5]:
```

	clientid	income	age	loan	default
count	2000.000000	2000.000000	1997.000000	2000.000000	2000.000000
mean	1000.500000	45331.600018	40.807559	4444.369695	0.141500
std	577.494589	14326.327119	13.624469	3045.410024	0.348624
min	1.000000	20014.489470	-52.423280	1.377630	0.000000
25%	500.750000	32796.459717	28.990415	1939.708847	0.000000
50%	1000.500000	45789.117313	41.317159	3974.719419	0.000000
75%	1500.250000	57791.281668	52.587040	6432.410625	0.000000
max	2000.000000	69995.685578	63.971796	13766.061239	1.000000

```
In [10]:
```

```
Out[10]:
```

	clientid	income	age	loan	default
422	423	69995.685578	52.719673	2084.370861	0

```
In [6]: base_credir[base_credir ['income']>= 69995.685578]
```

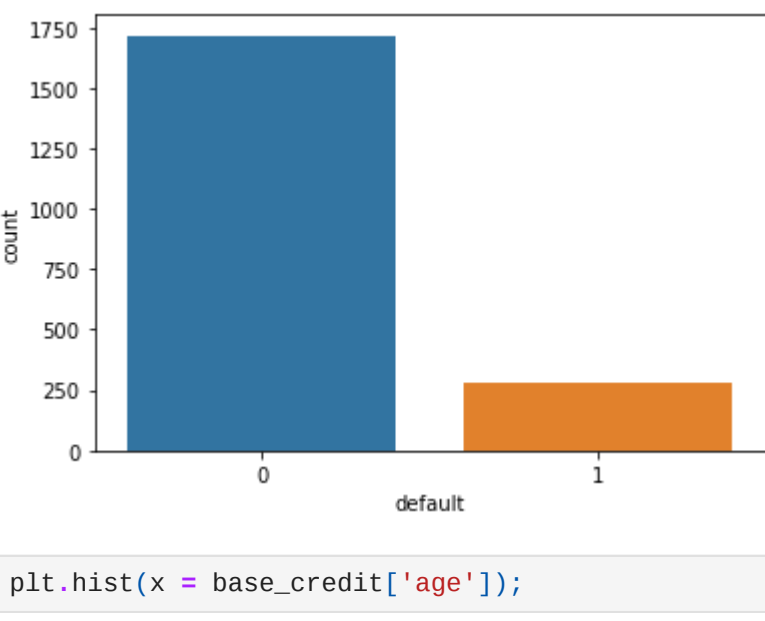
```
Out[6]:
```

	clientid	income	age	loan	default
422	423	69995.685578	52.719673	2084.370861	0

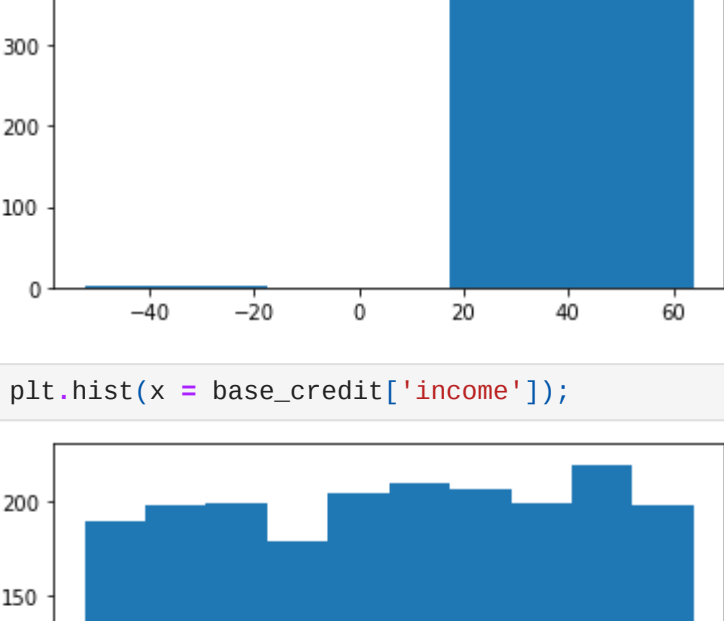
```
In [17]: np.unique(base_credir["default"],return_counts=True)
```

```
Out[17]: (array([0, 1], dtype=int64), array([1717, 283], dtype=int64))
```

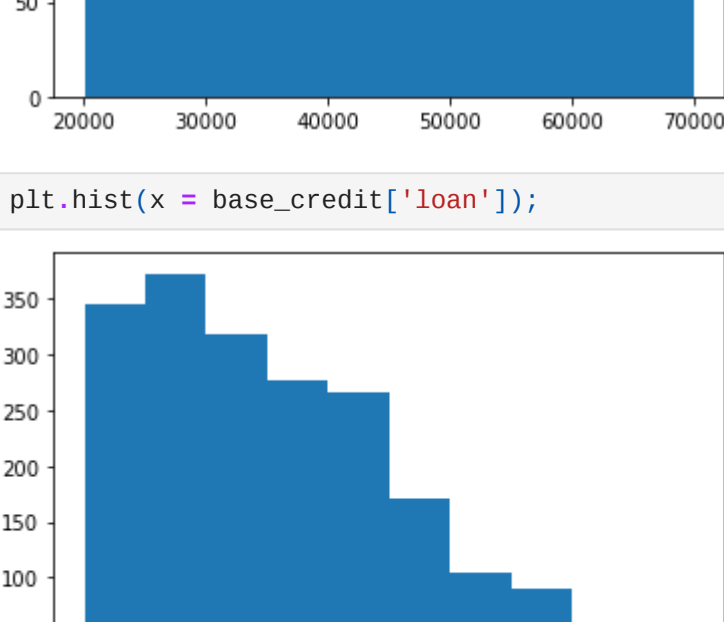
```
In [7]: sns.countplot(x = base_credir['default']);
```



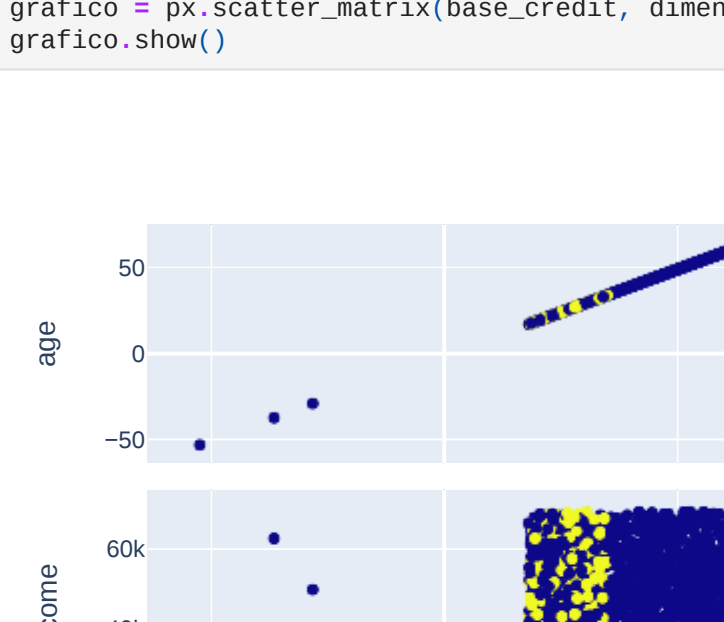
```
In [8]: plt.hist(x = base_credir['age']);
```



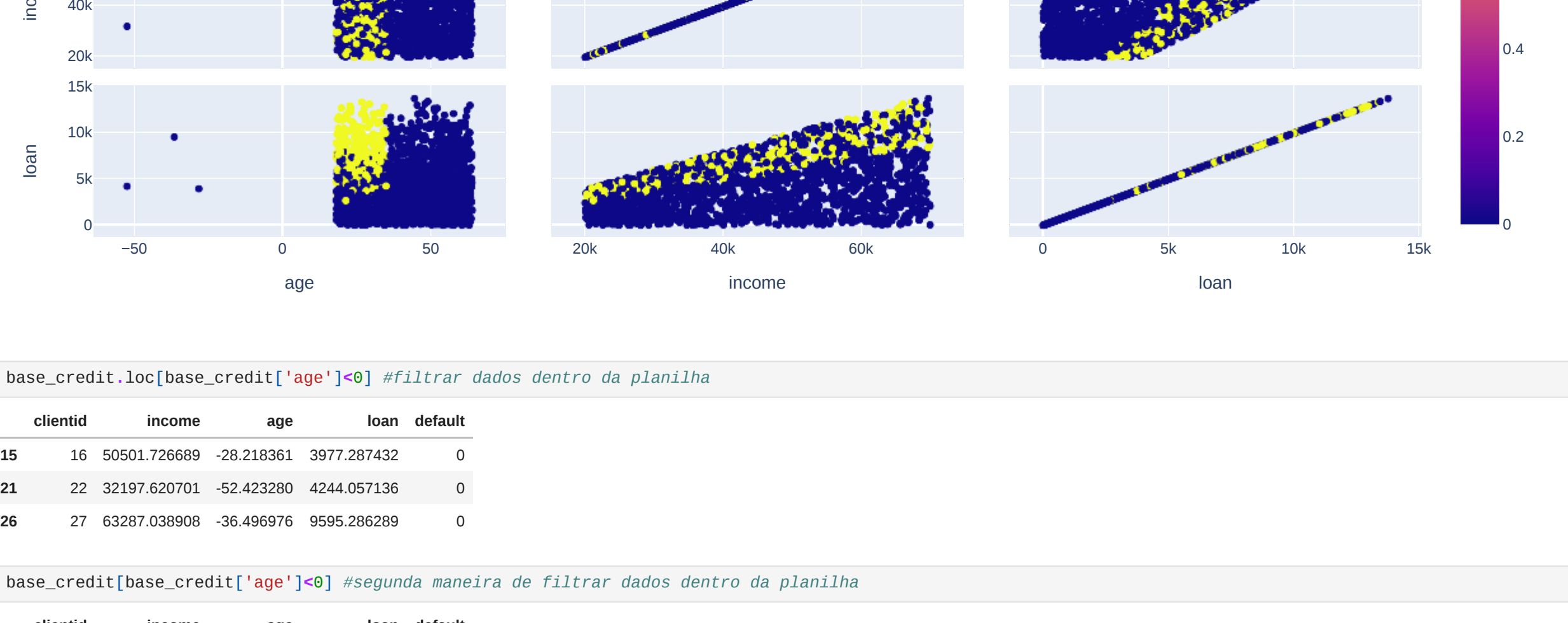
```
In [9]: plt.hist(x = base_credir['income']);
```



```
In [10]: plt.hist(x = base_credir['loan']);
```



```
In [11]: grafico = px.scatter_matrix(base_credir, dimensions=['age','income','loan'], color= 'default') #scatter = gráfico de dispersão 1= Não paga dívida, 0 =paga a dí
grafico.show()
```



```
In [34]: base_credir.loc[base_credir['age']<0] #filtrar dados dentro da planilha
```

```
Out[34]:
```

15	16	50501.726689	-28.218361	3977.287432	0
21	22	32197.620701	-52.423280	4244.057136	0
26	27	63287.038908	-36.496976	9595.286289	0

```
In [12]: base_credir[base_credir['age']<0] #segunda maneira de filtrar dados dentro da planilha
```

```
Out[12]:
```

	clientid	income	age	loan	default
15	16	50501.726689	-28.218361	3977.287432	0
21	22	32197.620701	-52.423280	4244.057136	0
26	27	63287.038908	-36.496976	9595.286289	0

```
In [13]: ##maneiras de tratar os valores de idade em negativo
#1 apagar as na linha inteira
base_credir2 = base_credir.drop('age',axis=1)#aqui estamos deletando a coluna 'age' axis=1 é coluna e x=0 é linha
base_credir2
```

```
Out[13]:
```

	clientid	income	loan	default
0	1	66155.925095	8106.532131	0
1	2	34415.153966	6564.745018	0
2	3	57317.170063	8020.953296	0
3	4	42709.534201	6103.642260	0
4	5	66952.688845	8770.099235	1
...
1995	1996	59221.044874	1926.729397	0
1996	1997	69516.127573	3503.176156	0
1997	1998	44311.449262	5522.786693	1
1998	1999	43756.056605	1622.722598	0
1999	2000	69436.579552	7378.833599	0

2000 rows × 4 columns

```
In [14]: #2 Apagar os registros com valores incoocistentes
base_credir3 = base_credir.drop(base_credir[base_credir['age']<0].index)
base_credir3
```

```
Out[14]:
```

	clientid	income	age	loan	default
0	1	66155.925095	59.017015	8106.532131	0
1	2	34415.153966	48.117153	6564.745018	0
2	3	57317.170063	63.108049	8020.953296	0
3	4	42709.534201	45.751972	6103.642260	0
4	5	66952.688845	18.584336	8770.099235	1
...
1995	1996	59221.044874	48.518179	1926.729397	0
1996	1997	69516.127573	23.162104	3503.176156	0
1997	1998	44311.449262	28.017167	5522.786693	1
1998	1999	43756.056605	63.971796	1622.722598	0
1999	2000	69436.579552	56.152617	7378.833599	0

1997 rows × 5 columns

```
In [15]: #3 preencher os valores manualmente(essa é a mais comum)
base_credir.mean()#estamos puxando a media de todas as colunas
```

```
Out[15]:
```

clientid	1000.500000
income	45331.600018
age	40.807559
loan	4444.369695
default	0.141500
dtype:	float64

```
In [16]: base_credir['age'].mean()#estamos acessando a média apenas da idade
```

```
Out[16]: 40.80755937848458
```

```
In [43]: #fazendo a média apenas dos valores da coluna age +0
base_credir['age'][base_credir['age']> 0].mean()
```

```
Out[43]: 40.92778044906149
```

```
In [17]: base_credir.loc[base_credir['age'] <0, 'age'] = 40.92
```

```
In [18]: base_credir.loc[base_credir['age']<0]
```

```
Out[18]:
```

	clientid	income	age	loan	default
--	----------	--------	-----	------	---------

```
In [53]: base_credir.head(27)
```

```
Out[53]:
```

0	1	66155.925095	59.017015	8106.532131	0
1	2	34415.153966	48.117153	6564.745018	0
2	3	57317.170063	63.108049	8020.953296	0
3	4	42709.534201	45.751972	6103.642260	0
4	5	66952.688845	18.584336	8770.099235	1
...
1995	1996	59221.044874	48.518179	1926.729397	0
1996	1997	69516.127573	23.162104	3503.176156	0
1997	1998	44311.449262	28.017167	5522.786693	1
1998	1999	43756.056605	63.971796	1622.722598	0
1999	2000	69436.579552	56.152617	7378.833599	0
...
28	29	59417.805406	NaN	2082.625938	0
30	31	48528.852796	NaN	6155.784670	0
31	32	23526.302555	NaN	2862.010139	0

```
In [19]: #Tratamento de dados faltantes
base_credir.isnull() #estamos verificando se ha celulas sem dados
```

```
Out[19]:
```

	clientid	income	age	loan	default
0	False	False	False	False	False
1	False	False	False	False	False
2	False	False	False	False	False
3	False	False	False	False	False
4	False	False	False	False	False
...
1995	False	False	False	False	False
1996	False	False	False	False	False
1997	False	False	False	False	False
1998	False	False	False	False	False
1999	False	False	False	False	False

2000 rows × 5 columns

```
In [63]: base_credir.isnull().sum() # aqui estamos somando a quantidade de celuas que estão sem dados em cada coluna
```

```
Out[63]:
```

clientid	0
income	0
age	3
loan	0
default	0
dtype:	int64

```
In [20]: base_credir.loc[pd.isnull(base_credir['age'])]
```

```
Out[20]:
```

	clientid	income	age	loan	default
28	29	59417.805406	NaN	2082.625938	0
30	31	48528.852796	NaN	6155.784670	0
31	32	23526.302555	NaN	2862.010139	0

```
In [66]: #aqui estamos alterando o valor das celulas em branco com o valor da media da idade fillna=preencherValornulo(ou seja ele vai preencher apenas em celulas que
base_credir['age'].fillna(base_credir['age'].mean(),inplace=True)
```

```
Out[67]:
```

	clientid	income	age	loan	default
--	----------	--------	-----	------	---------

```
In [21]: #estamos selecionando as celulas apenas para verificar que elas foram preenchidas 'isin' procura nas celas que estiverem dentro de parenteses
base_credir.loc[base_credir['clientid'].isin([29,31,32])]
```

```
Out[21]:
```

28	29	59417.805406	NaN	2082.625938	0
30	31	48528.852796	NaN	6155.784670	0
31	32	23526.302555	NaN	2862.010139	0

```
In [72]: # x normalmente são as variáveis previsoras e y são as classe
x_credit = base_credir.iloc[:,1:4].values # primeiro ':' incia que estamos selecionando todas as linhas, '1:4' estamos selcionando income, age e loan
```

```
Out[72]: array([[6.61559251e+04, 5.90170151e+01, 8.10653213e+03],
       [3.44151540e+04, 4.81171531e+01, 6.56474502e+03],
       [5.73171701e+04, 6.31080495e+01, 8.02095330e+03],
       ...,
       [4.43114493e+04, 2.80171609e+01, 5.52278669e+03],
       [4.37560566e+04, 6.39717958e+01, 1.62272260e+03],
       [6.94365796e+04, 5.61526170e+01, 7.37883360e+03]])
```

```
In [75]: y_credit = base_credir.iloc[:,4].values
x_credit
```

```
Out[75]: array([[6.61559251e+04, 5.90170151e+01, 8.10653213e+03],
       [3.44151540e+04, 4.81171531e+01, 6.56474502e+03],
       [5.73171701e+04, 6.31080495e+01, 8.02095330e+03],
       ...,
       [4.43114493e+04, 2.80171609e+01, 5.52278669e+03],
       [4.37560566e+04, 6.39717958e+01, 1.62272260e+03],
       [6.94365796e+04, 5.61526170e+01, 7.37883360e+03]])
```

```
In [84]: x_credit[:,0].min(), x_credit[:,0].max()# procurar pela pessoa com menor renda e maior renda
```

```
Out[84]: (20014.4894700497, 69995.6855783239)
```

```
In [86]: x_credit[:,0].max(), x_credit[:,1].max()
```

```
Out[86]: (69995.6855783239, 63.9719584112021)
```

```
In [89]: from sklearn.preprocessing import StandardScaler #aqui estamos padronizando os valores
scaler_credit = StandardScaler()
x_credit = scaler_credit.fit_transform(x_credit)
```

```
Out[89]: array([[ 1.45393393,  1.36538093,  1.20281942],
       [-0.76217555,  0.5426002 ,  0.69642695],
       [ 0.83682973,  1.67471189,  1.17471147],
       ...,
       [-0.07122592, -0.67448519,  0.35420081],
       [-0.11008289,  1.73936739, -0.92675625],
       [ 1.682906 ,  1.14917639,  0.96381030]])
```

```
In [ ]:
```