

```
import pandas as pd
import numpy as np
import matplotlib as plt
%matplotlib inline
```

```
from google.colab import drive
drive.mount('/content/drive')
```

Mounted at /content/drive

```
train = pd.read_csv("/content/drive/MyDrive/loan/loan-train.csv")
```

```
test = pd.read_csv("/content/drive/MyDrive/loan/loan-test.csv")
```

```
train.head()
```

	Loan_ID	Gender	Married	Dependents	Education	Self_Employed	ApplicantIncome	CoapplicantIncome	LoanAmount	Loan_Amount_Term	C
0	LP001002	Male	No	0	Graduate	No	5849	0.0	NaN	360.0	
1	LP001003	Male	Yes	1	Graduate	No	4583	1508.0	128.0	360.0	
2	LP001005	Male	Yes	0	Graduate	Yes	3000	0.0	66.0	360.0	
3	LP001006	Male	Yes	0	Not Graduate	No	2583	2358.0	120.0	360.0	
4	LP001008	Male	No	0	Graduate	No	6000	0.0	141.0	360.0	

```
train.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 614 entries, 0 to 613
Data columns (total 13 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Loan_ID               614 non-null   object
1   Gender                601 non-null   object
2   Married               611 non-null   object
3   Dependents            599 non-null   object
4   Education             614 non-null   object
5   Self_Employed         582 non-null   object
6   ApplicantIncome       614 non-null   int64
7   CoapplicantIncome     614 non-null   float64
8   LoanAmount            592 non-null   float64
9   Loan_Amount_Term      600 non-null   float64
10  Credit_History         564 non-null   float64
11  Property_Area         614 non-null   object
12  Loan_Status           614 non-null   object
dtypes: float64(4), int64(1), object(8)
memory usage: 62.5+ KB
```

```
train.shape
```

(614, 13)

```
train.describe()
```

index	ApplicantIncome	CoapplicantIncome	LoanAmount
count	614.0	614.0	592.0
mean	5403.459283387622	1621.2457980271008	146.41216216216216
std	6109.041673387178	2926.2483692241885	85.58732523570
min	150.0	0.0	0.0
25%	2877.5	0.0	100.0
50%	3812.5	1188.5	120.0
75%	5795.0	2297.25	160.0
max	81000.0	41667.0	700.0

Show 25 per page



Like what you see? Visit the [data table notebook](#) to learn more about interactive tables.

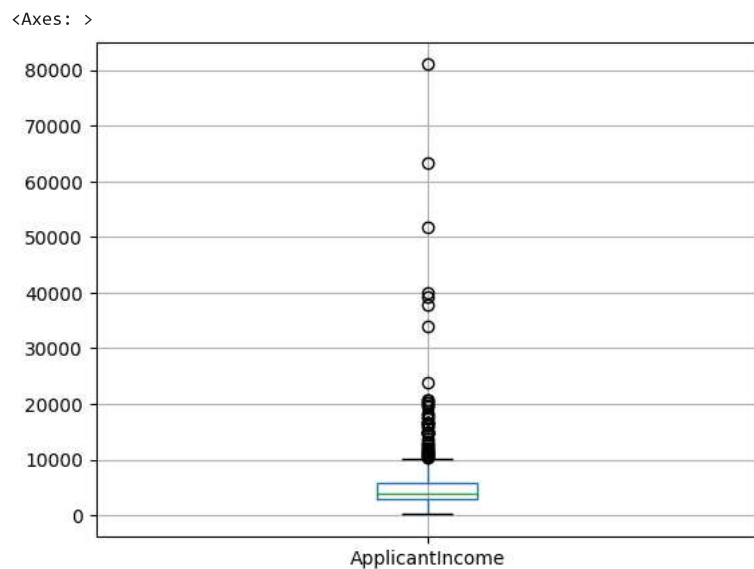
## ▼ To see how credit history affect the loan status

applicant credit history of 1 are more elligible for loan than one's who have credit history 0

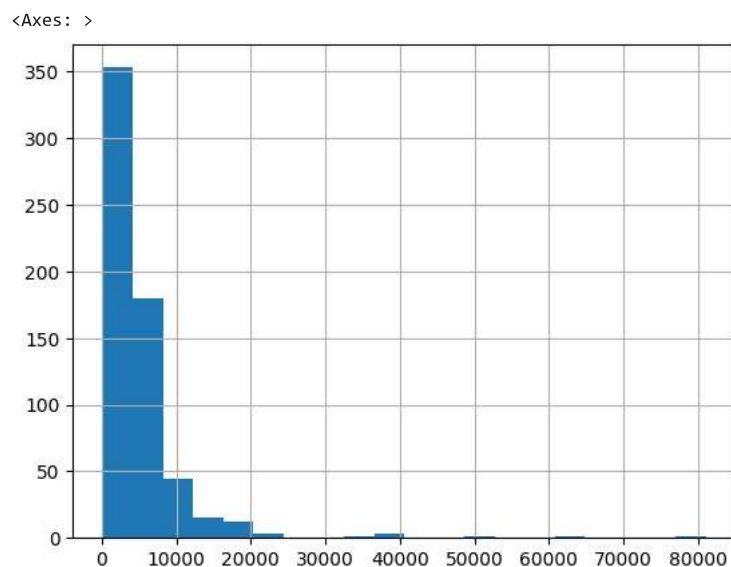
```
pd.crosstab(train['Credit_History'], train['Loan_Status'], margins=True)
```

Loan_Status	N	Y	All
Credit_History			
0.0	82	7	89
1.0	97	378	475
All	179	385	564

```
train.boxplot(column= 'ApplicantIncome')
```

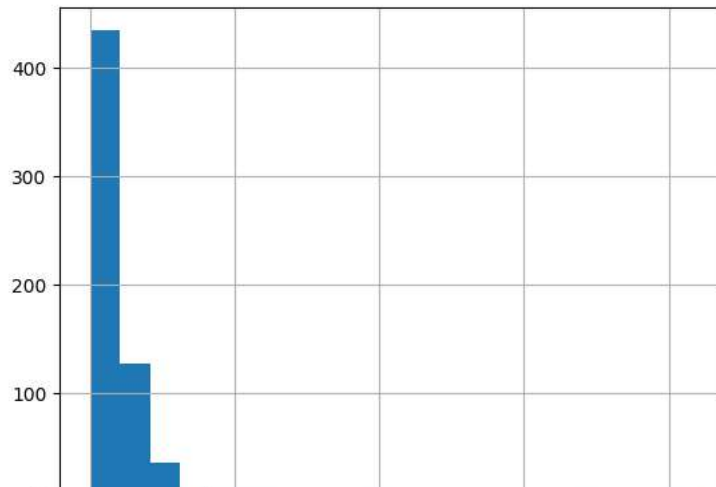


```
train['ApplicantIncome'].hist(bins=20)
```



```
train['CoapplicantIncome'].hist(bins=20)
```

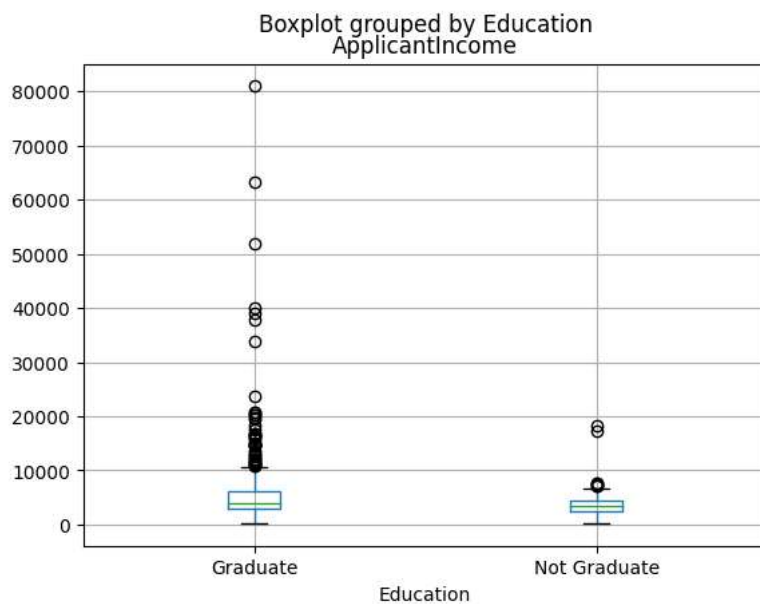
&lt;Axes: &gt;



Relation between Applicants Income and their education through boxplot

```
train.boxplot(column= 'ApplicantIncome', by= 'Education')
```

&lt;Axes: title={'center': 'ApplicantIncome'}, xlabel='Education'&gt;

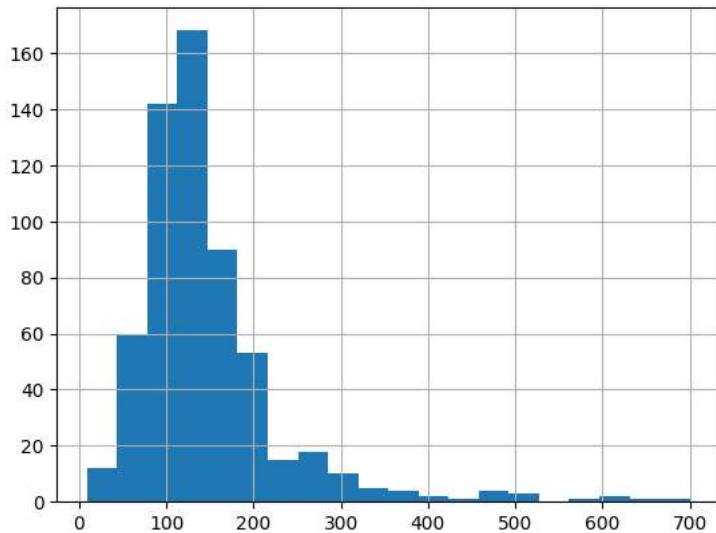


```
train.boxplot(column='LoanAmount')
```

&lt;Axes: &gt;

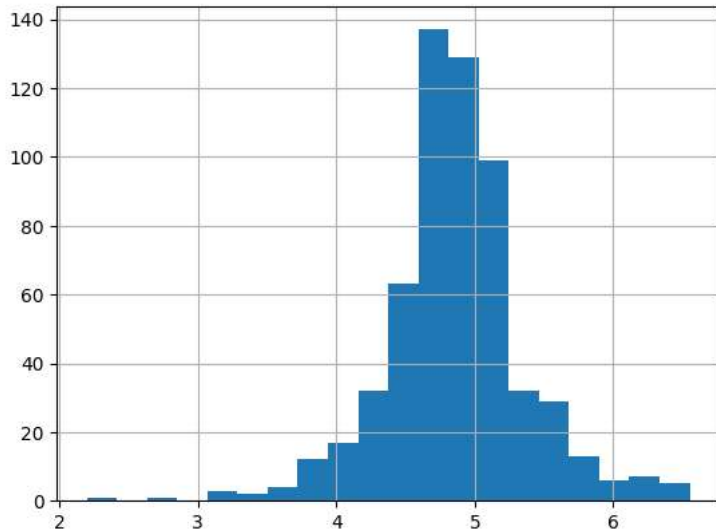
```
train['LoanAmount'].hist(bins=20)
```

&lt;Axes: &gt;



```
## Loan Amount is right skewed , for normalizing it apply log function
train['LoanAmount_log']=np.log(train['LoanAmount'])
train['LoanAmount_log'].hist(bins=20)
```

&lt;Axes: &gt;



```
train.isnull().sum()
```

```
Loan_ID          0
Gender           13
Married          3
Dependents       15
Education        0
Self_Employed   32
ApplicantIncome  0
CoapplicantIncome 0
LoanAmount       22
Loan_Amount_Term 14
Credit_History  50
Property_Area    0
Loan_Status      0
LoanAmount_log   22
dtype: int64
```

```
train['Gender'].fillna(train['Gender'].mode()[0], inplace=True)
```

```
train['Married'].fillna(train['Married'].mode()[0], inplace=True)
```

```
train['Dependents'].fillna(train['Dependents'].mode()[0], inplace=True)
```

```
train['Self_Employed'].fillna(train['Self_Employed'].mode()[0], inplace=True)
```

```
train.LoanAmount = train.LoanAmount.fillna(train.LoanAmount.mean())
train.LoanAmount_log = train.LoanAmount_log.fillna(train.LoanAmount_log.mean())
```

```
train['Loan_Amount_Term'].fillna(train['Loan_Amount_Term'].mode()[0], inplace=True)
```

```
train['Credit_History'].fillna(train['Credit_History'].mode()[0], inplace=True)
```

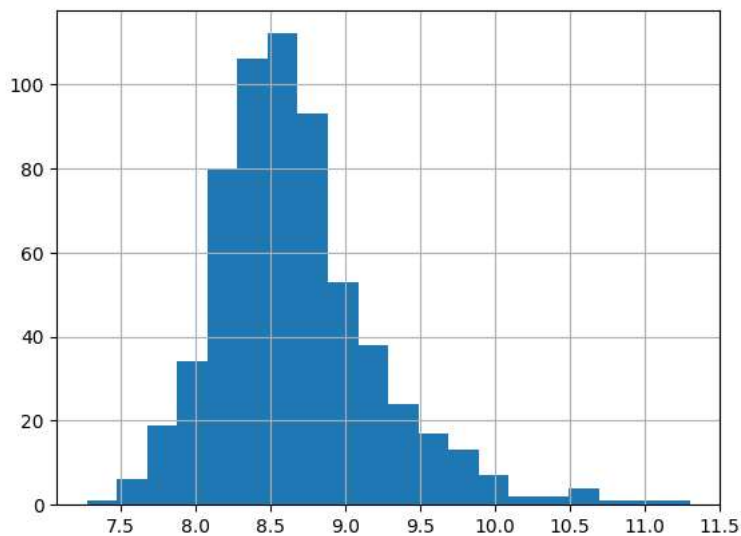
```
train.isnull().sum()
```

```
Loan_ID          0
Gender           0
Married          0
Dependents       0
Education        0
Self_Employed    0
ApplicantIncome  0
CoapplicantIncome 0
LoanAmount       0
Loan_Amount_Term 0
Credit_History   0
Property_Area    0
Loan_Status      0
LoanAmount_log    0
dtype: int64
```

```
train['TotalIncome']= train['ApplicantIncome'] + train['CoapplicantIncome']
train['TotalIncome_log']=np.log(train['TotalIncome'])
```

```
train['TotalIncome_log'].hist(bins=20)
```

<Axes: >



```
train.head()
```

6/10

[illegible]

```
for i in range(0, 5):
    X_test[:, i] = labelencoder_X.fit_transform(X_test[:, i])
```

```
X_test[:,7]= labelencoder_X.fit_transform(X_test[:,7])
```

```
labelencoder_y=LabelEncoder()
y_test= labelencoder_y.fit_transform(y_test)
```

y test

```
array([1, 0, 1, 0, 1, 0, 1, 1, 0, 1, 1, 1, 1, 0, 0, 1, 1, 0, 0, 1,
       1, 1, 1, 1, 0, 1, 0, 0, 1, 1, 1, 1, 0, 1, 1, 1, 0, 1, 1,
       1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 0, 1, 0, 1,
       1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 0, 1, 0, 0, 1, 0, 1, 1, 1,
       1, 1, 1, 0, 0, 0, 1, 0, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1,
       1, 0, 0, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1])
```

```
from sklearn.preprocessing import StandardScaler
ss=StandardScaler()
X_train=ss.fit_transform(X_train)
X_test=ss.fit_transform(X_test)
```

- Decision tree classifier

```
from sklearn.tree import DecisionTreeClassifier
DTClassifier= DecisionTreeClassifier(criterion='entropy', random_state=0)
DTClassifier.fit(X_train, y_train)
```

```
DecisionTreeClassifier
DecisionTreeClassifier(criterion='entropy', random_state=0)
```

```
y_pred= DTClassssifier.predict(X_test)
y_pred
```

```
array([0, 1, 0, 1, 1, 1, 0, 0, 0, 1, 1, 0, 0, 1, 1, 1, 0, 1, 1, 0, 0, 1,
       1, 1, 1, 1, 1, 1, 0, 0, 0, 1, 0, 1, 1, 0, 1, 1, 0, 1, 1, 0, 0, 1,
       1, 1, 1, 0, 0, 0, 1, 1, 1, 1, 0, 0, 0, 1, 1, 0, 0, 1, 0, 1, 1,
       1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 0, 1, 0, 1, 1, 0, 1, 1, 1, 1,
       1, 0, 0, 1, 0, 0, 1, 0, 1, 1, 0, 1, 1, 0, 1, 1, 1, 0, 1, 1, 1,
       1, 0, 1, 1, 0, 0, 1, 1, 0, 1, 0, 0, 1])
```

```
from sklearn import metrics
print('the accuracy of decision tree is : ', metrics.accuracy_score(y_pred,y_test))
```

the accuracy of decision tree is : 0.7073170731707317

```
from sklearn.naive_bayes import GaussianNB
NBClassifier = GaussianNB()
NBClassifier.fit(X_train,y_train)
```

- ▼ GaussianNB

```
GaussianNB()
```

```
y_pred= NBClassifier.predict(X_test)
```

y\_pred

```
array([1, 1, 1, 1, 1, 0, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 1,
       1, 1, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1, 1,
```

```
1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 0, 1, 1,
1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1,
1, 1, 1, 1, 1, 0, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 1, 0, 1])
```

```
print('The accuracy of Naive Bayes is :', metrics.accuracy_score(y_pred,y_test))
```

```
The accuracy of Naive Bayes is : 0.8292682926829268
```

```
test.head()
```

	Loan_ID	Gender	Married	Dependents	Education	Self_Employed	ApplicantIncome	Coa
0	LP001015	Male	Yes	0	Graduate	No	5720	
1	LP001022	Male	Yes	1	Graduate	No	3076	
2	LP001031	Male	Yes	2	Graduate	No	5000	
3	LP001035	Male	Yes	2	Graduate	No	2340	
4	LP001051	Male	No	0	Not	No	3076	

```
test.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 367 entries, 0 to 366
Data columns (total 12 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Loan_ID                367 non-null   object
1   Gender                 356 non-null   object
2   Married                367 non-null   object
3   Dependents             357 non-null   object
4   Education              367 non-null   object
5   Self_Employed          344 non-null   object
6   ApplicantIncome        367 non-null   int64
7   CoapplicantIncome      367 non-null   int64
8   LoanAmount             362 non-null   float64
9   Loan_Amount_Term       361 non-null   float64
10  Credit_History         338 non-null   float64
11  Property_Area          367 non-null   object
dtypes: float64(3), int64(2), object(7)
memory usage: 34.5+ KB
```

```
test.isnull().sum()
```

```
Loan_ID                0
Gender                 11
Married                0
Dependents             10
Education              0
Self_Employed          23
ApplicantIncome        0
CoapplicantIncome      0
LoanAmount             5
Loan_Amount_Term       6
Credit_History        29
Property_Area          0
dtype: int64
```

```
test['Gender'].fillna(test['Gender'].mode()[0], inplace=True)
test['Dependents'].fillna(test['Dependents'].mode()[0], inplace=True)
test['Self_Employed'].fillna(test['Self_Employed'].mode()[0], inplace=True)
test['Loan_Amount_Term'].fillna(test['Loan_Amount_Term'].mode()[0], inplace=True)
test['Credit_History'].fillna(test['Credit_History'].mode()[0], inplace=True)
```

```
test.isnull().sum()
```

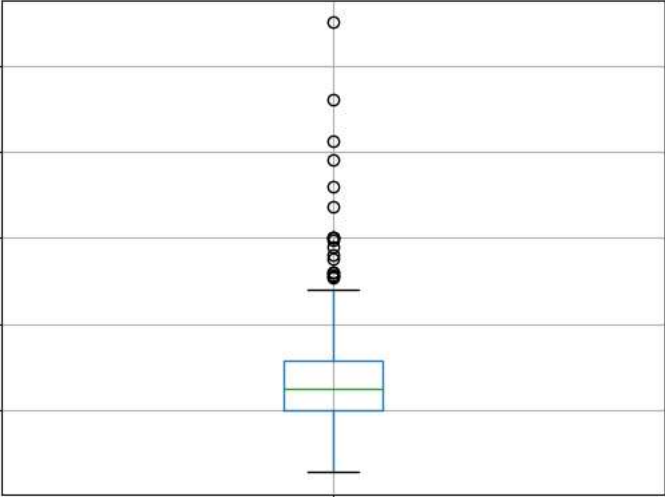
```
Loan_ID                0
Gender                 0
Married                0
Dependents             0
Education              0
Self_Employed          0
ApplicantIncome        0
CoapplicantIncome      0
LoanAmount             5
```



```
Loan_Amount_Term    0
Credit_History      0
Property_Area        0
dtype: int64

test.boxplot(column='LoanAmount')
```

<Axes: >



A boxplot showing the distribution of 'LoanAmount'. The y-axis ranges from 100 to 500. The box is blue with a green median line at approximately 125. Whiskers extend from 40 to 240. Numerous outliers are plotted as open circles, ranging from approximately 250 to 550.

LoanAmount

```
test.LoanAmount= test.LoanAmount.fillna(test.LoanAmount.mean())
```

```
test['LoanAmount_log']=np.log(test['LoanAmount'])
```

```
test.isnull().sum()
```

```
Loan_ID            0
Gender             0
Married            0
Dependents         0
Education          0
Self_Employed     0
ApplicantIncome    0
CoapplicantIncome  0
LoanAmount         0
Loan_Amount_Term   0
Credit_History     0
Property_Area      0
LoanAmount_log     0
dtype: int64
```

```
test['TotalIncome']= test['ApplicantIncome']+test['CoapplicantIncome']
test['TotalIncome_log']= np.log(test['TotalIncome'])
```

```
test.head()
```

	Loan_ID	Gender	Married	Dependents	Education	Self_Employed	ApplicantIncome	Coa
0	LP001015	Male	Yes	0	Graduate	No	5720	
1	LP001022	Male	Yes	1	Graduate	No	3076	
2	LP001031	Male	Yes	2	Graduate	No	5000	
3	LP001035	Male	Yes	2	Graduate	No	2340	
4	LP001051	Male	No	0	Not Graduate	No	3276	

```
test_variable= test.iloc[:,np.r_[1:5,9:11,13:15]].values
```

```
array([1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1,
1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1,
1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 0, 1, 1, 1, 0, 1, 1, 1,
0, 0, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 1, 0, 1, 1, 1, 1,
1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 0, 1, 1, 1,
1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 1, 1, 0, 0, 1, 0, 1, 1, 1, 1, 1, 1,
1, 1, 1, 1, 1, 1, 1, 0, 1, 0, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 0,
1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 0, 0, 1, 0, 1, 1, 1, 1, 0, 0, 1,
1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 0, 1,
0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1,
1, 1, 1, 1, 0, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0,
1, 0, 1, 0, 1, 1, 1, 1, 0, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1,
1, 1, 0, 1, 0, 1, 1, 1, 1, 0, 0, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1,
1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1,
1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1,
1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1,
1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1,
1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1])
```