

In [82]:

```

1 from sklearn.ensemble import GradientBoostingClassifier
2 import numpy as np
3 import pandas as pd
4 from sklearn.preprocessing import StandardScaler
5 from sklearn.model_selection import train_test_split
6 from sklearn.metrics import accuracy_score, confusion_matrix
7 from sklearn import preprocessing
8 import matplotlib.pyplot as plt
9 import seaborn as sns
10 import warnings
11 warnings.filterwarnings("ignore")

```

In [69]:

```

1 df = pd.read_csv("C:/Users/HP/Downloads/income_evaluation.csv")
2 df.columns

```

Out[69]:

```

Index(['age', 'workclass', 'fnlwgt', 'education', 'education-num',
      'marital-status', 'occupation', 'relationship', 'race', 'sex',
      'capital-gain', 'capital-loss', 'hours-per-week', 'native-country',
      'income'],
      dtype='object')

```

In [70]:

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

Out[70]:

	age	workclass	fnlwgt	education	education-num	marital-status	occupation	relationship	race
0	39	State-gov	77516	Bachelors	13	Never-married	Adm-clerical	Not-in-family	White
1	50	Self-emp-not-inc	83311	Bachelors	13	Married-civ-spouse	Exec-managerial	Husband	White
2	38	Private	215646	HS-grad	9	Divorced	Handlers-cleaners	Not-in-family	White
3	53	Private	234721	11th	7	Married-civ-spouse	Handlers-cleaners	Husband	Black
4	28	Private	338409	Bachelors	13	Married-civ-spouse	Prof-specialty	Wife	Black

In [71]:

```
1 df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 32561 entries, 0 to 32560
Data columns (total 15 columns):
 #   Column                Non-Null Count  Dtype
---  -
 0   age                   32561 non-null  int64
 1   workclass              32561 non-null  object
 2   fnlwgt                 32561 non-null  int64
 3   education              32561 non-null  object
 4   education-num          32561 non-null  int64
 5   marital-status         32561 non-null  object
 6   occupation              32561 non-null  object
 7   relationship           32561 non-null  object
 8   race                   32561 non-null  object
 9   sex                    32561 non-null  object
10   capital-gain           32561 non-null  int64
11   capital-loss           32561 non-null  int64
12   hours-per-week         32561 non-null  int64
13   native-country         32561 non-null  object
14   income                 32561 non-null  object
dtypes: int64(6), object(9)
memory usage: 3.7+ MB
```

```
1 # Data Preprocessing
```

In [74]:

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

Out[74]:

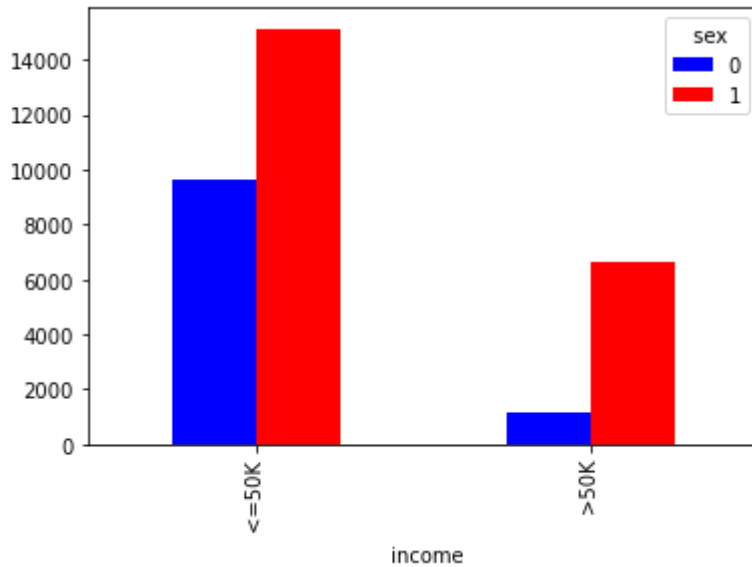
	age	workclass	fnlwgt	education	education-num	marital-status	occupation	relationship	race
0	39	State-gov	77516	Bachelors	13	Never-married	Adm-clerical	Not-in-family	White
1	50	Self-emp-not-inc	83311	Bachelors	13	Married-civ-spouse	Exec-managerial	Husband	White
2	38	Private	215646	HS-grad	9	Divorced	Handlers-cleaners	Not-in-family	White
3	53	Private	234721	11th	7	Married-civ-spouse	Handlers-cleaners	Husband	Black
4	28	Private	338409	Bachelors	13	Married-civ-spouse	Prof-specialty	Wife	Black

In [75]:

```
1 # plot of sex vs income
2 pd.crosstab(df[' income'],df[' sex']).plot(kind='bar',color=['blue','red'])
```

Out[75]:

<AxesSubplot:xlabel=' income'>



In [76]:

```
1 np.unique(df[' marital-status'])
```

Out[76]:

```
array([' Divorced', ' Married-AF-spouse', ' Married-civ-spouse',
       ' Married-spouse-absent', ' Never-married', ' Separated',
       ' Widowed'], dtype=object)
```

In [77]:

```
1 np.unique(df[' workclass'])
```

Out[77]:

```
array([' ?', ' Federal-gov', ' Local-gov', ' Never-worked', ' Private',
       ' Self-emp-inc', ' Self-emp-not-inc', ' State-gov', ' Without-pa
y'],
      dtype=object)
```

In [78]:

```
1 df.drop(columns=' fnlwgt',inplace=True)
```

In [79]:

```
1 column_names = ['age', 'workclass', 'education', 'education_num', 'marital_status',  
2                 'race', 'sex', 'capital_gain', 'capital_loss', 'hours_per_week', 'native_country',  
3                 'income']  
4 df.columns = column_names  
5 df.columns
```

Out[79]:

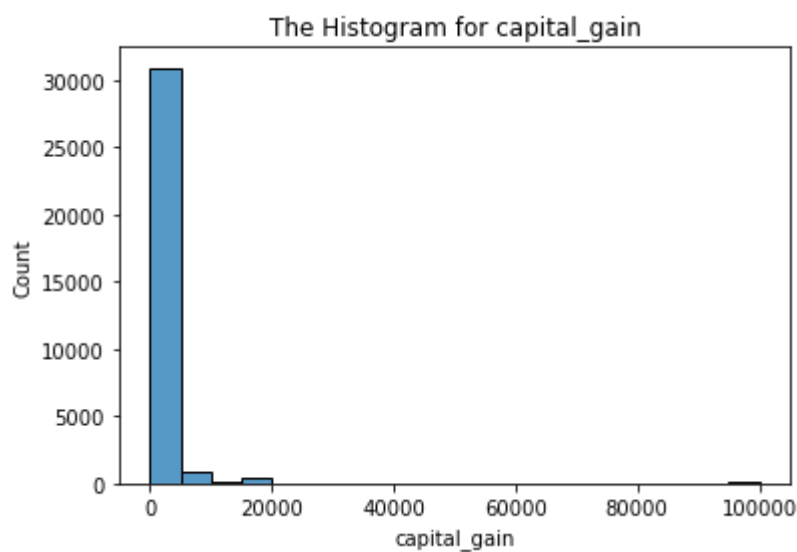
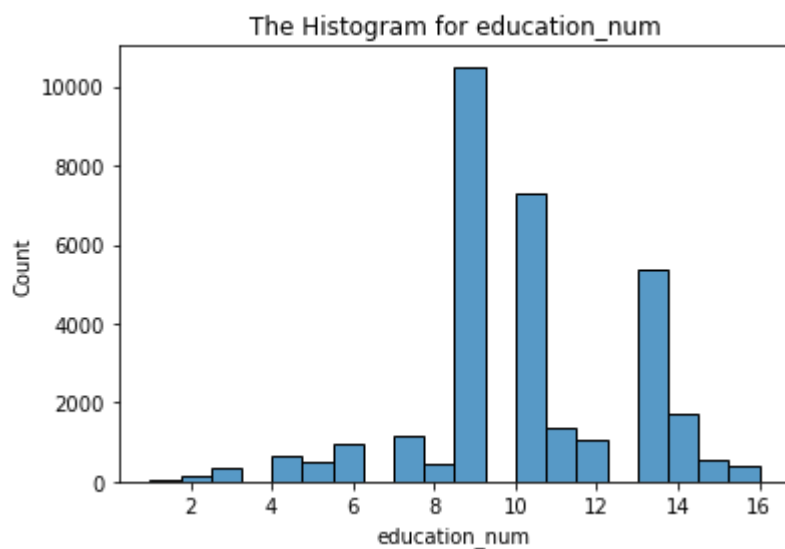
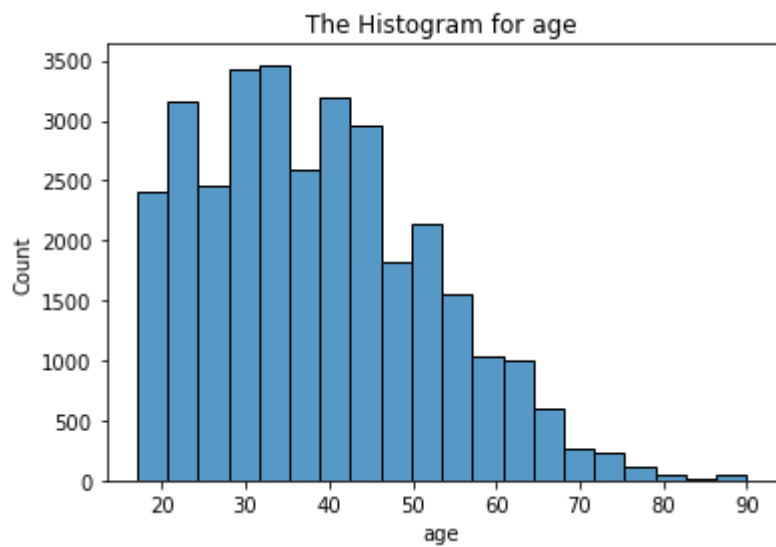
```
Index(['age', 'workclass', 'education', 'education_num', 'marital_status',  
      'occupation', 'relationship', 'race', 'sex', 'capital_gain',  
      'capital_loss', 'hours_per_week', 'native_country', 'income'],  
      dtype='object')
```

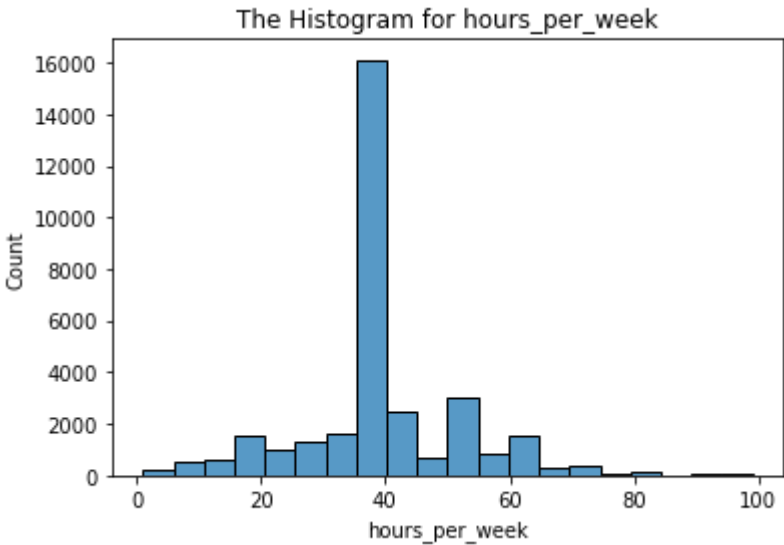
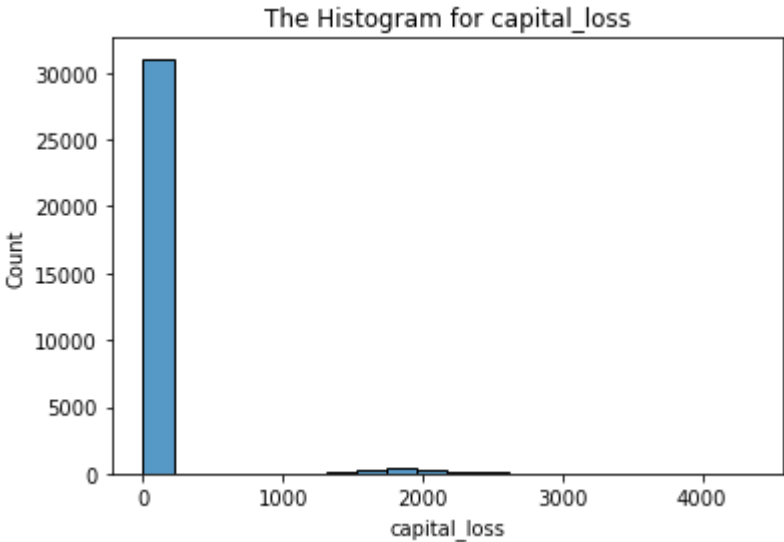
In [80]:

```
1 categorical = [var for var in df.columns if df[var].dtype == 'o' ]  
2 numerical = [num for num in df.columns if df[num].dtype == 'int64']
```

In [83]:

```
1 for i in numerical:
2     sns.histplot(x=df[i], palette='Set1',bins=20)
3     plt.title("The Histogram for {}".format(i))
4     plt.show()
```



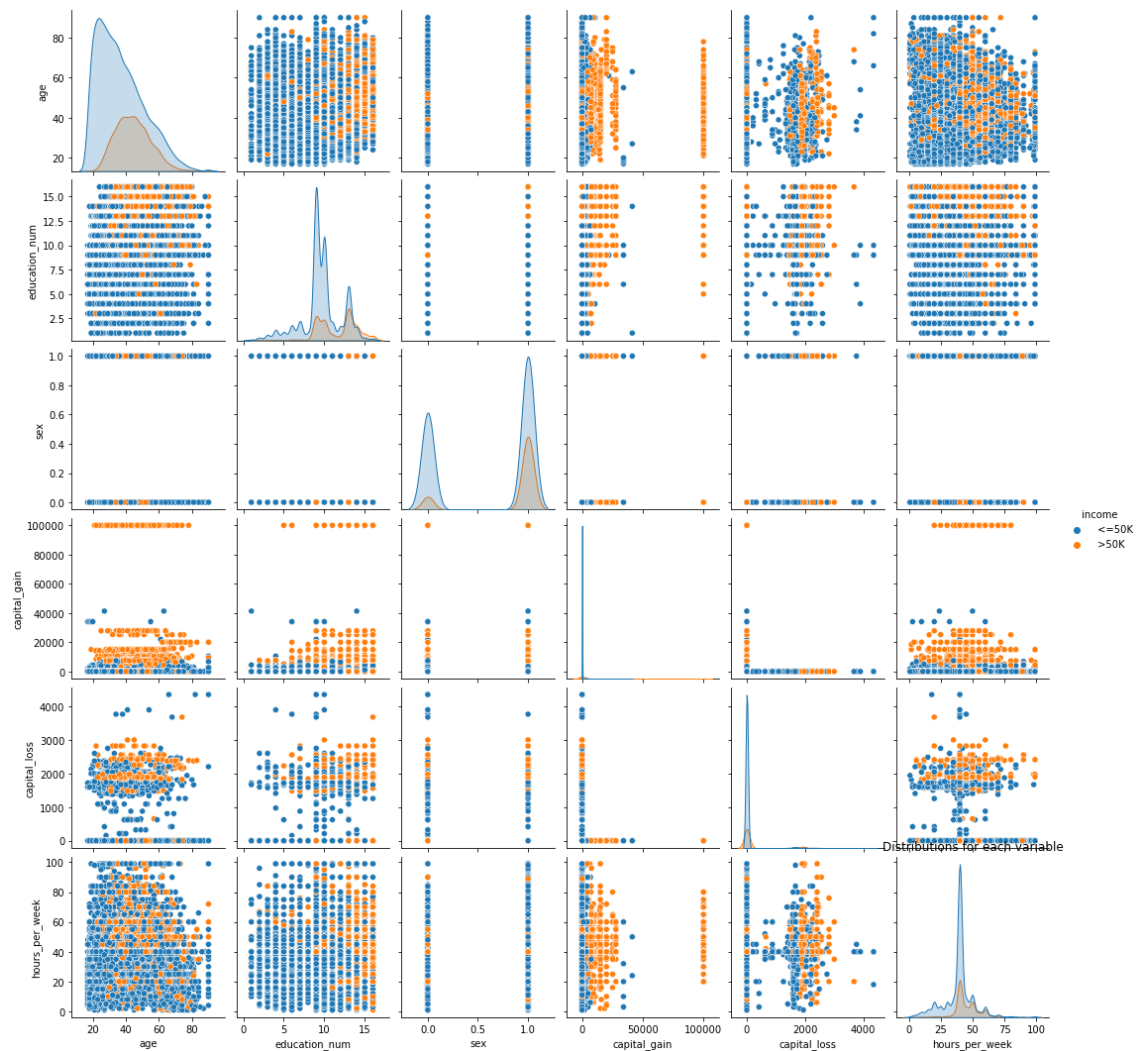


In [93]:

```

1 sns.pairplot(data=df, hue="income")
2 plt.title('Distributions for each variable')
3 plt.show()

```



In [101]:

```

1 def label_encoder(a):
2     le = LabelEncoder()
3     df[a] = le.fit_transform(df[a])
4
5 label_list = ['workclass', 'education', 'marital_status',
6               'occupation', 'relationship', 'race', 'sex', 'native_country', 'income']
7
8 for i in label_list:
9     label_encoder(i)

```

In [104]:

```

1 x = df.drop('income', axis=1)
2 y = df['income']

```

In [105]:

```
1 X_train, X_test, y_train, y_test = train_test_split( x,y, test_size=0.2)
```

In [110]:

```
1 # Define Gradient Boosting Classifier with hyperparameters
2 GBA=GradientBoostingClassifier(n_estimators=500,learning_rate=0.05,random_state=100)
3 GBA.fit(X_train,y_train)
4 print("Test Score:",GBA.score(X_test,y_test))
```

Test Score: 0.8694917856594503

In [111]:

```
1 print(confusion_matrix(y_test, GBA.predict(X_test)))
```

```
[[4644  249]
 [ 601 1019]]
```

Hyperparameter tuning

In [118]:

```
1 from sklearn.model_selection import cross_val_score
2 from sklearn.model_selection import GridSearchCV
3
4 grid = {
5     'learning_rate':[0.01,0.05,0.1],
6     'n_estimators':np.arange(100,500,100),
7 }
8
9 gb = GradientBoostingClassifier()
10 gb_cv = GridSearchCV(gb, grid, cv = 4)
11 gb_cv.fit(X_train,y_train)
12 print("Best Parameters:",gb_cv.best_params_)
13 print("Train Score:",gb_cv.best_score_)
14 print("Test Score:",gb_cv.score(X_test,y_test))
```

Best Parameters: {'learning_rate': 0.1, 'n_estimators': 300}
Train Score: 0.8707002457002456
Test Score: 0.8727161062490404