

Adult Census Income Prediction Project

INCOME PREDICTION



Problem Statement

- The Goal is to predict whether a person has an income of more than 50K a year or not.
- This is basically a binary classification problem where a person is classified into the >50K group or <=50K group.

Import Libraries

In [1]:

```
1 import warnings
2 warnings.filterwarnings('ignore')
3 import numpy as np
4 import pandas as pd
5 import seaborn as sns
6 import missingno as mno
7 import matplotlib.pyplot as plt
8 from sklearn.preprocessing import LabelEncoder
9 from imblearn.over_sampling import RandomOverSampler
10 from sklearn.linear_model import LogisticRegression
11 from sklearn.neighbors import KNeighborsClassifier
12 from sklearn.svm import SVC
13 from sklearn.naive_bayes import GaussianNB
14 from sklearn.tree import DecisionTreeClassifier
15 from sklearn.ensemble import RandomForestClassifier
16 from xgboost import XGBClassifier
17 from sklearn.ensemble import AdaBoostClassifier
18 from sklearn.metrics import accuracy_score
19 from sklearn.metrics import f1_score
20 from sklearn.metrics import confusion_matrix
21 from sklearn.metrics import classification_report
```

In [2]:

```
1 df = pd.read_csv("adult.csv", na_values="?", skipinitialspace = True)
2 df.head(5)
```

Out[2]:

	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

Exploratory Data Analysis (EDA)

In [3]:

```
1 df.shape
```

Out[3]:

(32561, 15)

In [4]:

```
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        30725 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       30718 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  country          31978 non-null   object  
 14  salary            32561 non-null   object  
dtypes: int64(6), object(9)
memory usage: 3.7+ MB
```

In [5]:

```
1 df.columns
```

Out[5]:

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

In [6]:

```
1 df.dtypes
```

Out[6]:

```
age           int64
workclass     object
fnlwgt        int64
education     object
education-num int64
marital-status object
occupation    object
relationship   object
race          object
sex            object
capital-gain  int64
capital-loss  int64
hours-per-week int64
country       object
salary         object
dtype: object
```

In [7]:

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

Out[7]:

	count	mean	std	min	25%	50%	75%
age	32561.0	38.581647	13.640433	17.0	28.0	37.0	48.0
fnlwgt	32561.0	189778.366512	105549.977697	12285.0	117827.0	178356.0	237051.0
education-num	32561.0	10.080679	2.572720	1.0	9.0	10.0	12.0
capital-gain	32561.0	1077.648844	7385.292085	0.0	0.0	0.0	0.0
capital-loss	32561.0	87.303830	402.960219	0.0	0.0	0.0	0.0
hours-per-week	32561.0	40.437456	12.347429	1.0	40.0	40.0	45.0



In [8]:

```
1 df.isnull().sum()
```

Out[8]:

```
age          0
workclass    1836
fnlwgt       0
education    0
education-num 0
marital-status 0
occupation   1843
relationship  0
race         0
sex          0
capital-gain 0
capital-loss 0
hours-per-week 0
country      583
salary        0
dtype: int64
```

In [9]:

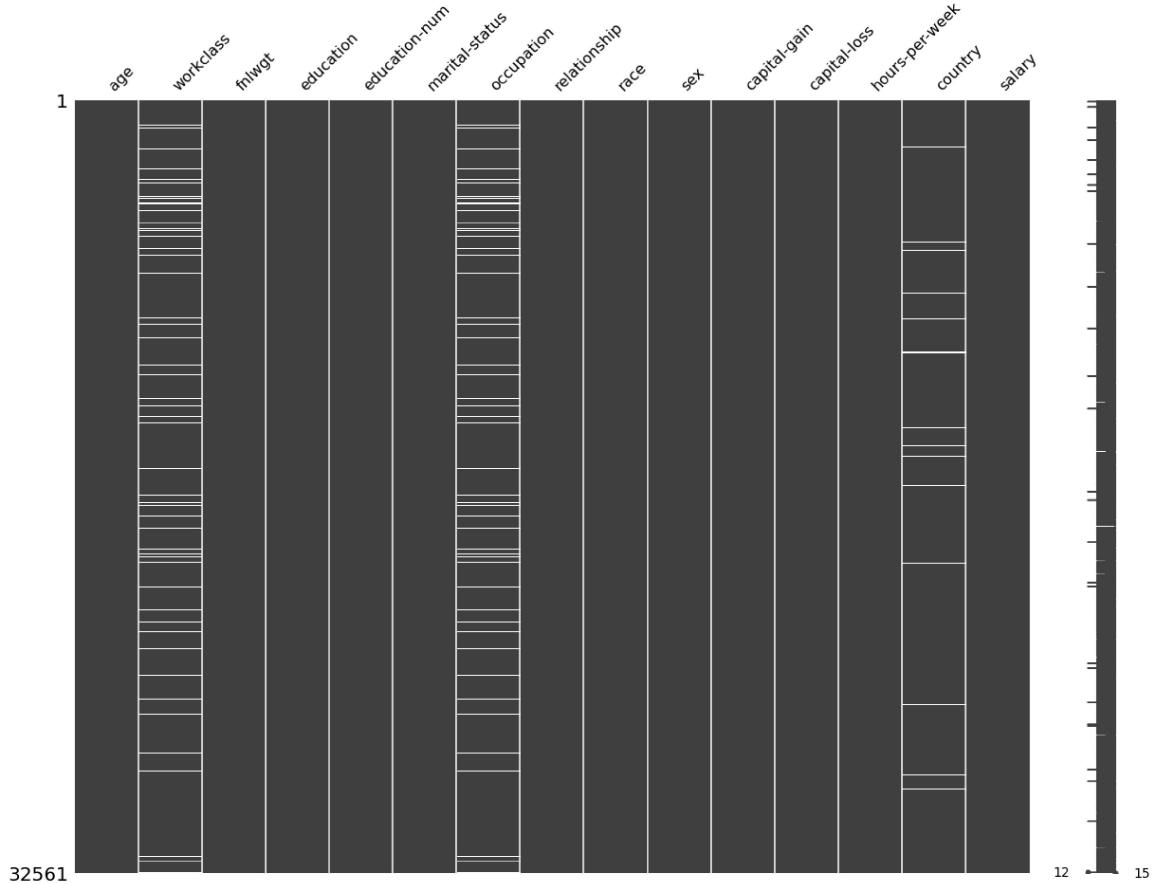
```
1 round(df.isnull().sum() / df.shape[0] * 100, 2).astype(str) + "%"
```

Out[9]:

```
age          0.0%
workclass    5.64%
fnlwgt       0.0%
education    0.0%
education-num 0.0%
marital-status 0.0%
occupation   5.66%
relationship  0.0%
race         0.0%
sex          0.0%
capital-gain 0.0%
capital-loss 0.0%
hours-per-week 0.0%
country      1.79%
salary        0.0%
dtype: object
```

In [10]:

```
1 mno.matrix(df, figsize = (20,15))  
2 plt.show()
```



In [11]:

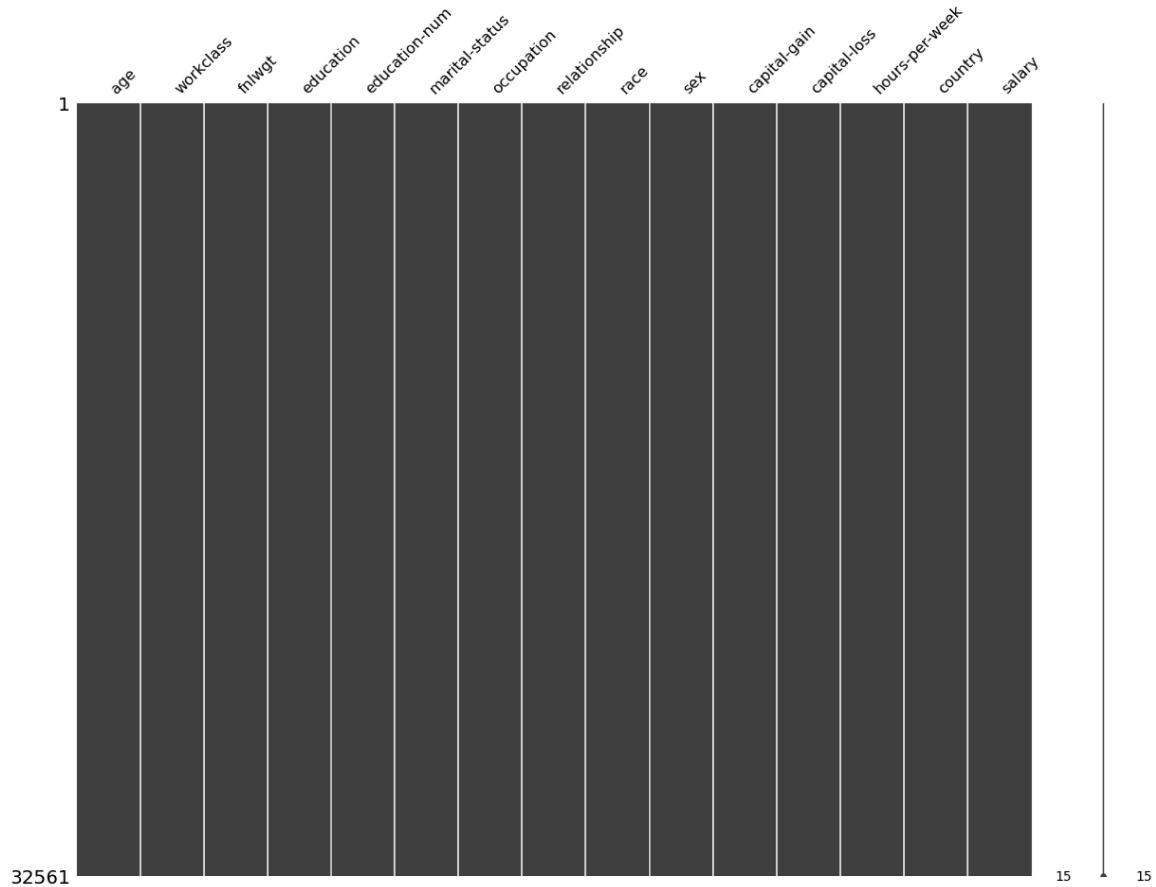
```
1 nan_values_columns = ['workclass', 'occupation', 'country']
```

In [12]:

```
1 for col in nan_values_columns:  
2     df[col].fillna(df[col].mode()[0], inplace=True)
```

In [13]:

```
1 mno.matrix(df, figsize = (20,15))  
2 plt.show()
```



In [14]:

```
1 df.unique()
```

Out[14]:

```
age              73  
workclass        8  
fnlwgt      21648  
education       16  
education-num    16  
marital-status     7  
occupation       14  
relationship      6  
race              5  
sex                2  
capital-gain     119  
capital-loss      92  
hours-per-week    94  
country            41  
salary              2  
dtype: int64
```

In [15]:

```
1 round(df.nunique() / df.shape[0] * 100, 2).astype(str) + "%"
```

Out[15]:

```
age           0.22%
workclass     0.02%
fnlwgt        66.48%
education     0.05%
education-num 0.05%
marital-status 0.02%
occupation    0.04%
relationship   0.02%
race          0.02%
sex            0.01%
capital-gain  0.37%
capital-loss  0.28%
hours-per-week 0.29%
country        0.13%
salary         0.01%
dtype: object
```

In [16]:

```
1 df.duplicated().sum()
```

Out[16]:

24

In [17]:

```
1 df = df.drop_duplicates()
2 df.head(5)
```

Out[17]:

	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 [18]:

```
1 df.shape
```

Out[18]:

(32537, 15)

In [19]:

```
1 salary = df['salary'].value_counts(normalize=True)
2 round(salary * 100, 2).astype('str') + '%'
```

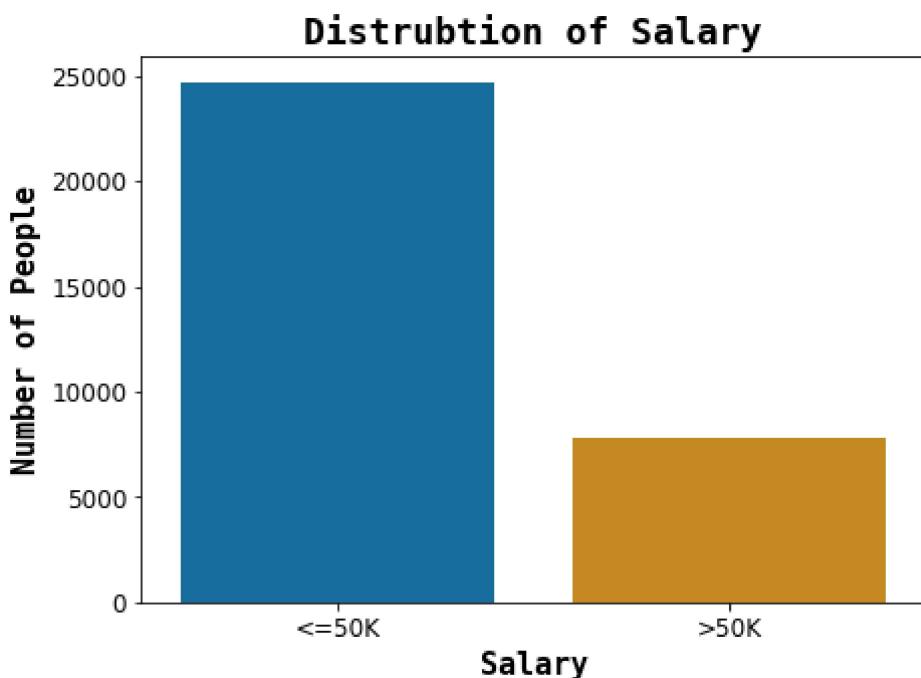
Out[19]:

```
<=50K    75.91 %
>50K    24.09 %
Name: salary, dtype: object
```

Data Visualization

In [20]:

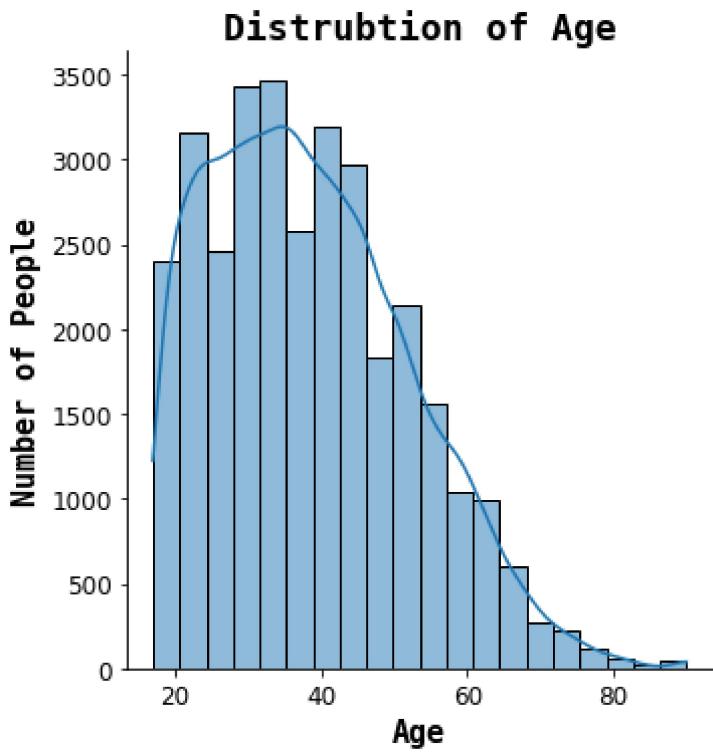
```
1 # barplot plot for 'Age'
2 salary= df['salary'].value_counts()
3
4 plt.figure(figsize=(7,5))
5 plt.tick_params(labelsize=12)
6 sns.barplot(salary.index, salary.values, palette='colorblind')
7 plt.title('Distrubtion of Salary', fontdict={
8     'fontname':'Monospace', 'fontsize':18, 'fontweight':'bold'})
9 plt.xlabel('Salary', fontdict={
10    'fontname':'Monospace', 'fontsize':15, 'fontweight':'bold'})
11 plt.ylabel('Number of People', fontdict={
12     'fontname':'Monospace', 'fontsize':15, 'fontweight':'bold'})
13 plt.show()
```



In [21]:

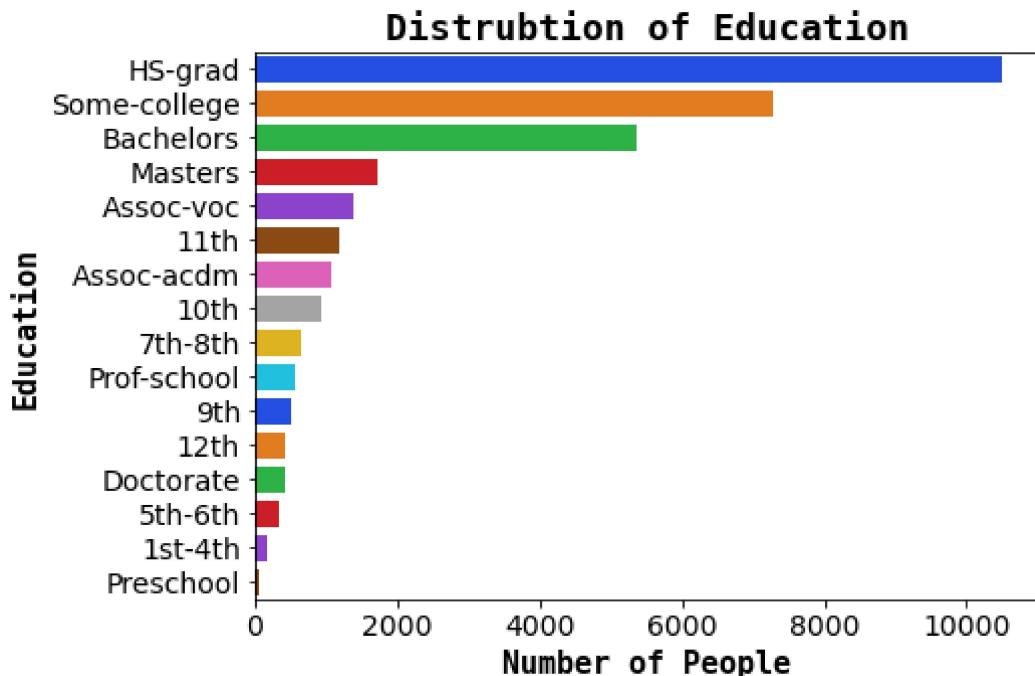
```
1 # distribution plot for 'Age'
2 age = df['age'].value_counts()
3
4 plt.figure(figsize=(20,15))
5 sns.displot(df['age'], bins=20, kde=True)
6 plt.title('Distrubtion of Age', fontdict={
7     'fontname':'Monospace', 'fontsize':18, 'fontweight':'bold'})
8 plt.xlabel('Age', fontdict={
9     'fontname':'Monospace', 'fontsize':15, 'fontweight':'bold'})
10 plt.ylabel('Number of People', fontdict={
11     'fontname':'Monospace', 'fontsize':15, 'fontweight':'bold'})
12 plt.tick_params(labelsize=12)
13 plt.show()
```

<Figure size 1440x1080 with 0 Axes>



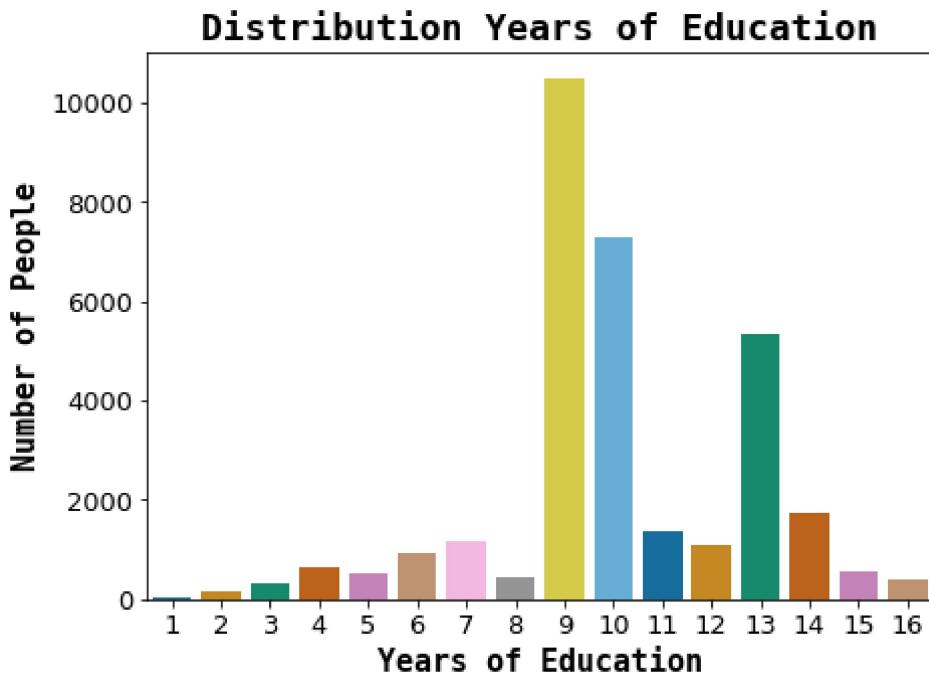
In [22]:

```
1 # Barplot for "Education"
2 education = df['education'].value_counts()
3
4 plt.figure(figsize=(7,5))
5 sns.barplot(education.values, education.index, palette='bright')
6 plt.title('Distribution of Education', fontdict={
7     'fontname':'Monospace', 'fontsize':18, 'fontweight':'bold'})
8 plt.xlabel("Number of People", fontdict={
9     'fontname':'Monospace', 'fontsize':15, 'fontweight':'bold'})
10 plt.ylabel("Education", fontdict={
11     'fontname':'Monospace', 'fontsize':15, 'fontweight':'bold'})
12 plt.tick_params(labelsize=14)
13 plt.show()
```



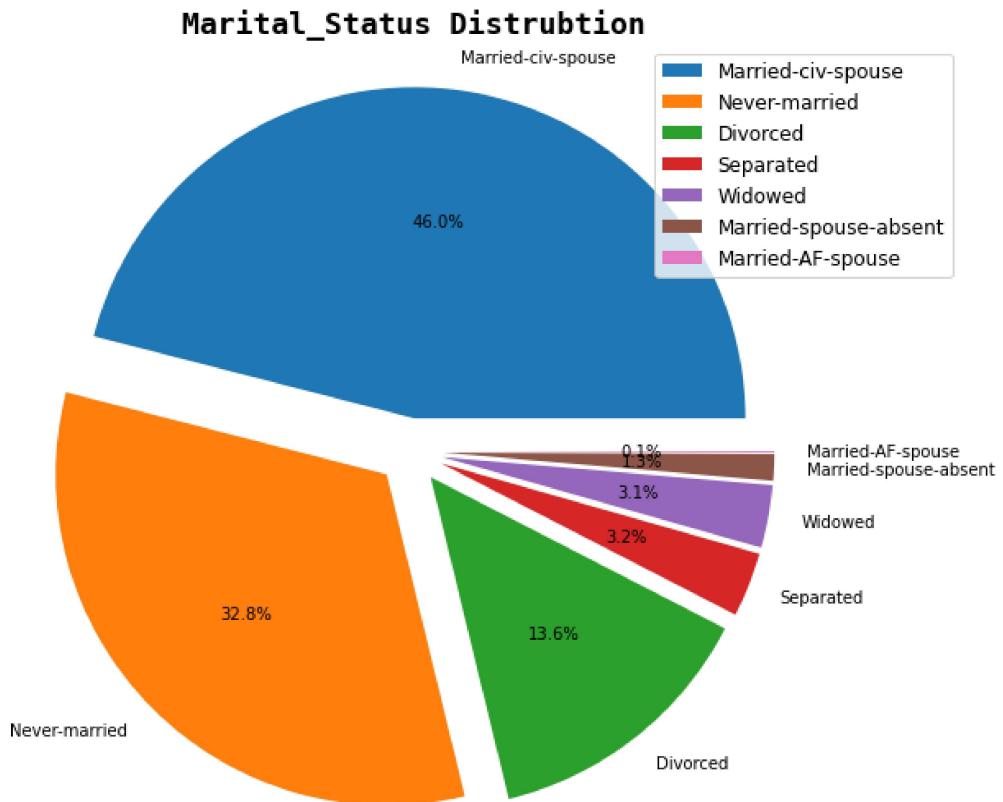
In [23]:

```
1 # Barplot for "Years of Education"
2
3 education_num = df['education-num'].value_counts()
4 plt.figure(figsize=(7,5))
5 sns.barplot(education_num.index, education_num.values, palette='colorblind')
6 plt.title('Distribution Years of Education', fontdict={
7     'fontname': 'Monospace', 'fontsize': 18, 'fontweight': 'bold'})
8 plt.xlabel('Years of Education', fontdict={
9     'fontname': 'Monospace', 'fontsize': 15, 'fontweight': 'bold'})
10 plt.ylabel('Number of People', fontdict={
11     'fontname': 'Monospace', 'fontsize': 15, 'fontweight': 'bold'})
12 plt.tick_params(labelsize=13)
13 plt.show()
```



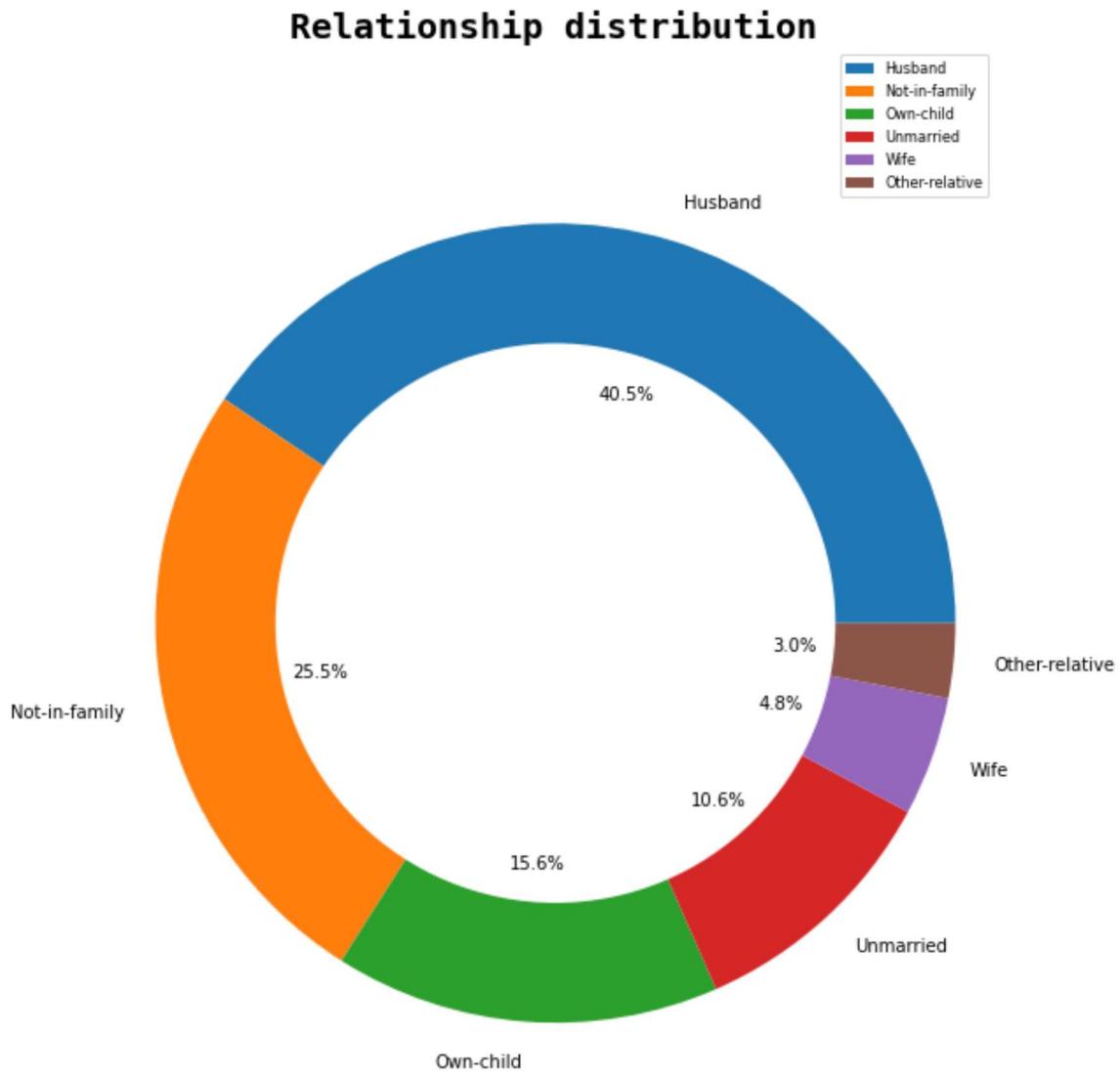
In [24]:

```
1 # Pie chart for "Marital_Status"
2 marital = df['marital-status'].value_counts()
3
4 plt.figure(figsize=(12,9))
5 plt.pie(marital.values, labels=marital.index,
6         explode=(0.10,0.10,0.10,0.10,0.10,0.10,0.10), autopct='%1.1f%%')
7 plt.title("Marital_Status Distrubtion",fontdict={
8     'fontname': 'Monospace', 'fontsize': 18, 'fontweight': 'bold'})
9 plt.legend()
10 plt.legend(prop={'size':12})
11 plt.axis('equal')
12 plt.show()
```



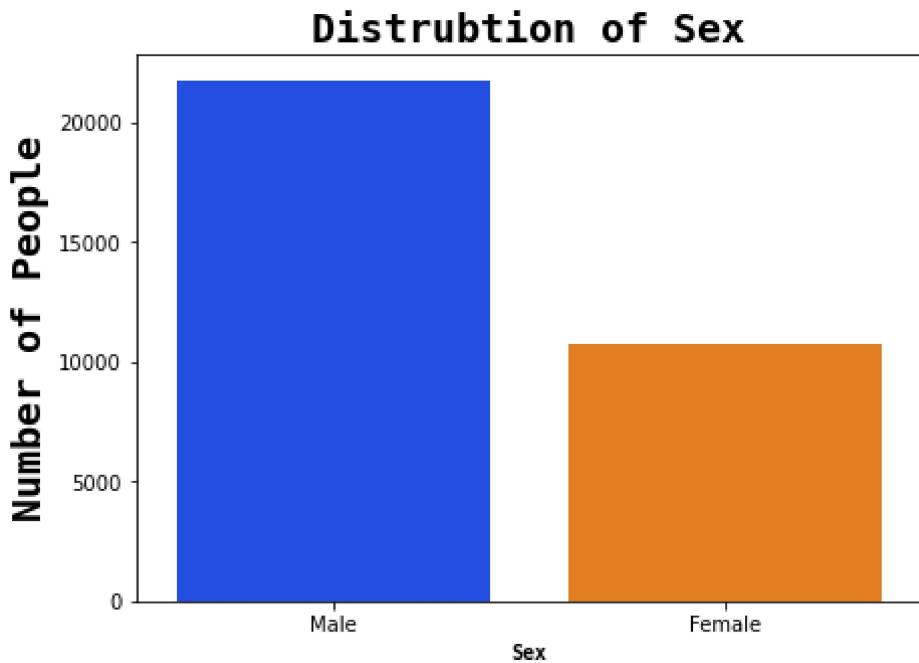
In [25]:

```
1 # Chart for 'Age'
2 relationship = df['relationship'].value_counts()
3
4 plt.figure(figsize=(9,12))
5 plt.pie(relationship.values,
6         labels=relationship.index,
7         autopct='%.1f%%')
8 center_part_circle = plt.Circle((0,0), 0.7, fc='White')
9 plt.title('Relationship distribution', fontdict={
10     'fontname': 'Monospace', 'fontsize': 20, 'fontweight': 'bold'})
11 fig = plt.gcf()
12 fig.gca().add_artist(center_part_circle)
13 plt.axis('equal')
14 plt.legend(prop={'size':8})
15 plt.show()
```



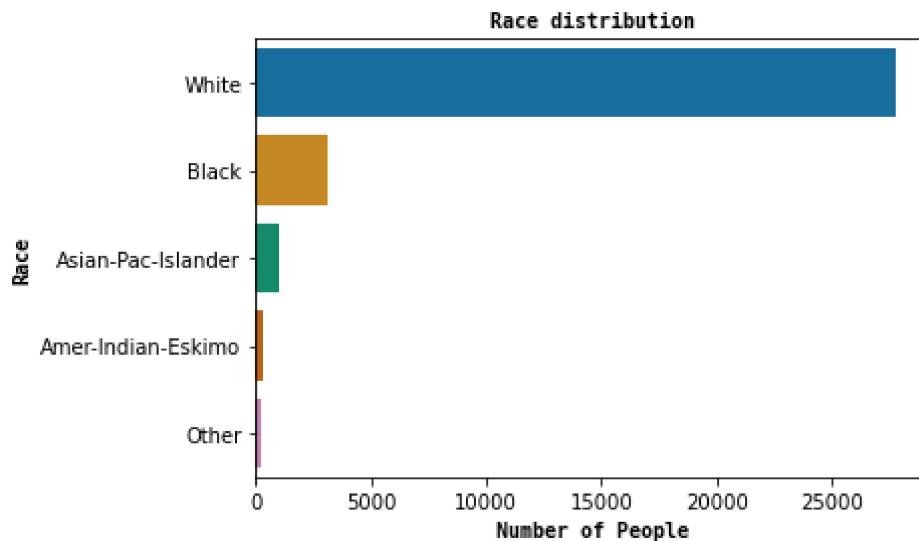
In [26]:

```
1 # barplot for 'Sex'  
2  
3 sex = df['sex'].value_counts()  
4 plt.figure(figsize=(7,5))  
5 sns.barplot(sex.index, sex.values, palette='bright')  
6 plt.title('Distribution of Sex', fontdict={  
7     'fontname': 'Monospace', 'fontsize': 20, 'fontweight': 'bold'})  
8 plt.xlabel('Sex', fontdict={  
9     'fontname': 'Monospace', 'fontsize': 10, 'fontweight': 'bold'})  
10 plt.ylabel('Number of People', fontdict={  
11     'fontname': 'Monospace', 'fontsize': 20, 'fontweight': 'bold'})  
12 plt.tick_params(labelsize=10)  
13 plt.show()
```



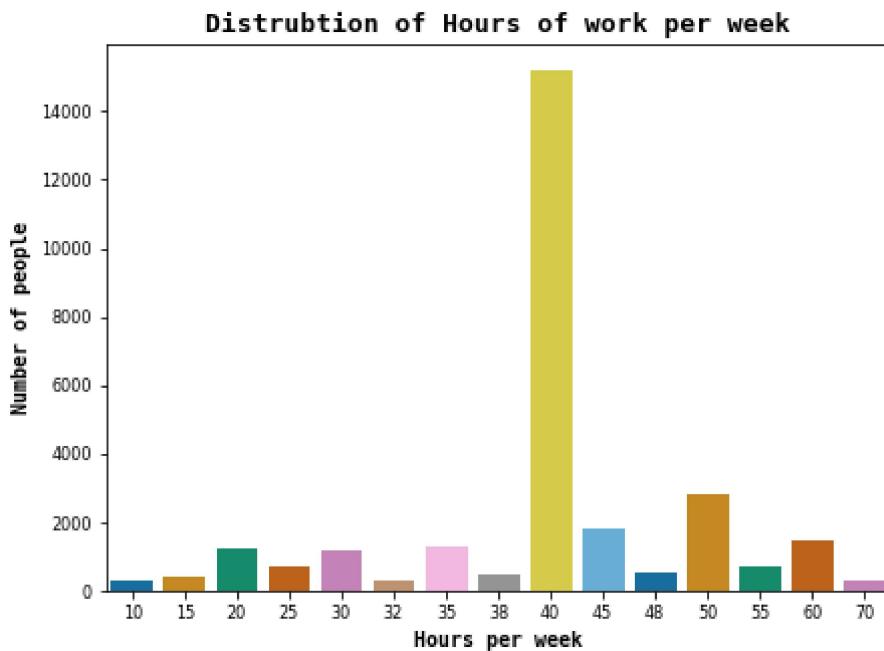
In [27]:

```
1 # barplot for 'Race'
2 race = df['race'].value_counts()
3
4 plt.figure(figsize=(6,4))
5 sns.barplot(race.values, race.index, palette='colorblind')
6 plt.title('Race distribution', fontdict={
7     'fontname': 'Monospace', 'fontsize': 10, 'fontweight': 'bold'})
8 plt.xlabel("Number of People", fontdict={
9     'fontname': 'Monospace', 'fontsize':10, 'fontweight':'bold'})
10 plt.ylabel("Race", fontdict={
11     'fontname':'Monospace', 'fontsize':10, 'fontweight':'bold'})
12 plt.tick_params(labelsize=10)
13 plt.show()
```



In [28]:

```
1 # barplot for 'Hours per week'
2 hour = df['hours-per-week'].value_counts().head(15)
3
4 plt.figure(figsize=(7,5))
5 sns.barplot(hour.index, hour.values, palette='colorblind')
6 plt.title('Distribution of Hours of work per week', fontdict={
7     'fontname':'Monospace', 'fontsize':13, 'fontweight':'bold'})
8 plt.xlabel('Hours per week', fontdict={
9     'fontname':'Monospace', 'fontsize':10, 'fontweight':'bold'})
10 plt.ylabel('Number of people', fontdict={
11     'fontname':'Monospace', 'fontsize':10, 'fontweight':'bold'})
12 plt.tick_params(labelsize=8)
13 plt.show()
```

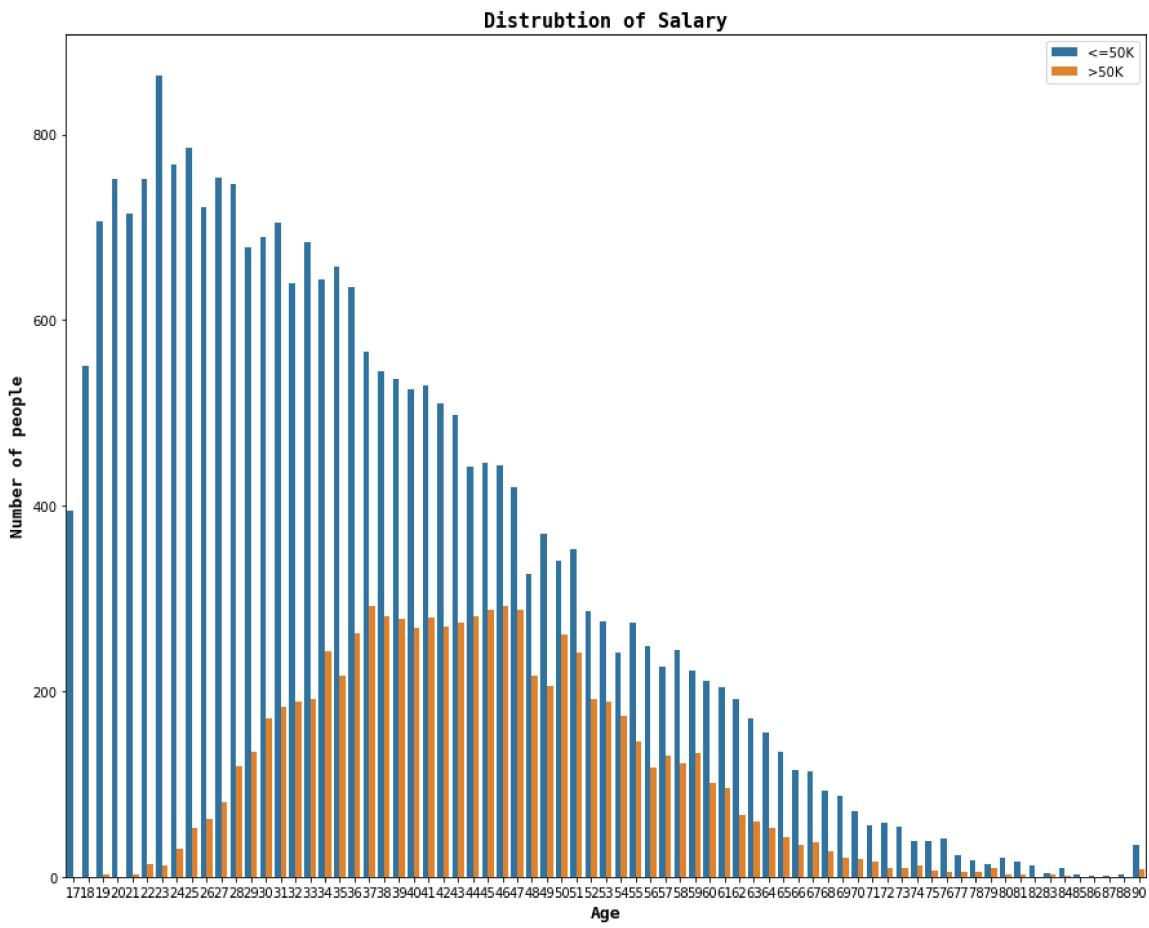


Bivariate Analysis

Bivariate analysis is a statistical technique used to analyze the relationship between two variables. In coding, bivariate analysis can be performed using a variety of programming languages and statistical software packages.

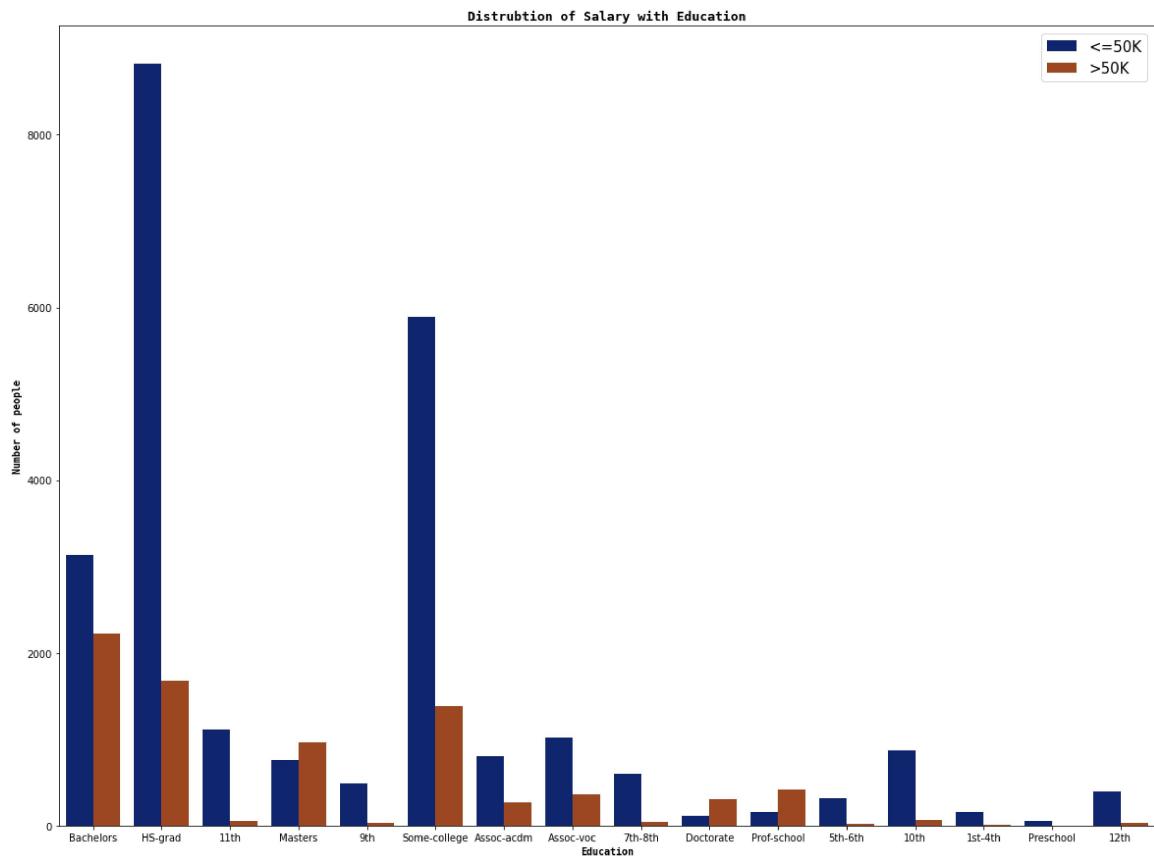
In [29]:

```
1 # Countplot of Salary with Age
2 plt.figure(figsize=(15,12))
3 sns.countplot(df['age'], hue=df['salary'])
4 plt.title('Distribution of Salary', fontdict={
5     'fontname':'Monospace', 'fontsize':15, 'fontweight':'bold'})
6 plt.xlabel('Age', fontdict={
7     'fontname':'Monospace', 'fontsize':13, 'fontweight':'bold'})
8 plt.ylabel('Number of people', fontdict={
9     'fontname':'Monospace', 'fontsize':13, 'fontweight':'bold'})
10 plt.legend(loc=1, prop={'size':10})
11 plt.tick_params(labelsize=10)
12 plt.show()
```



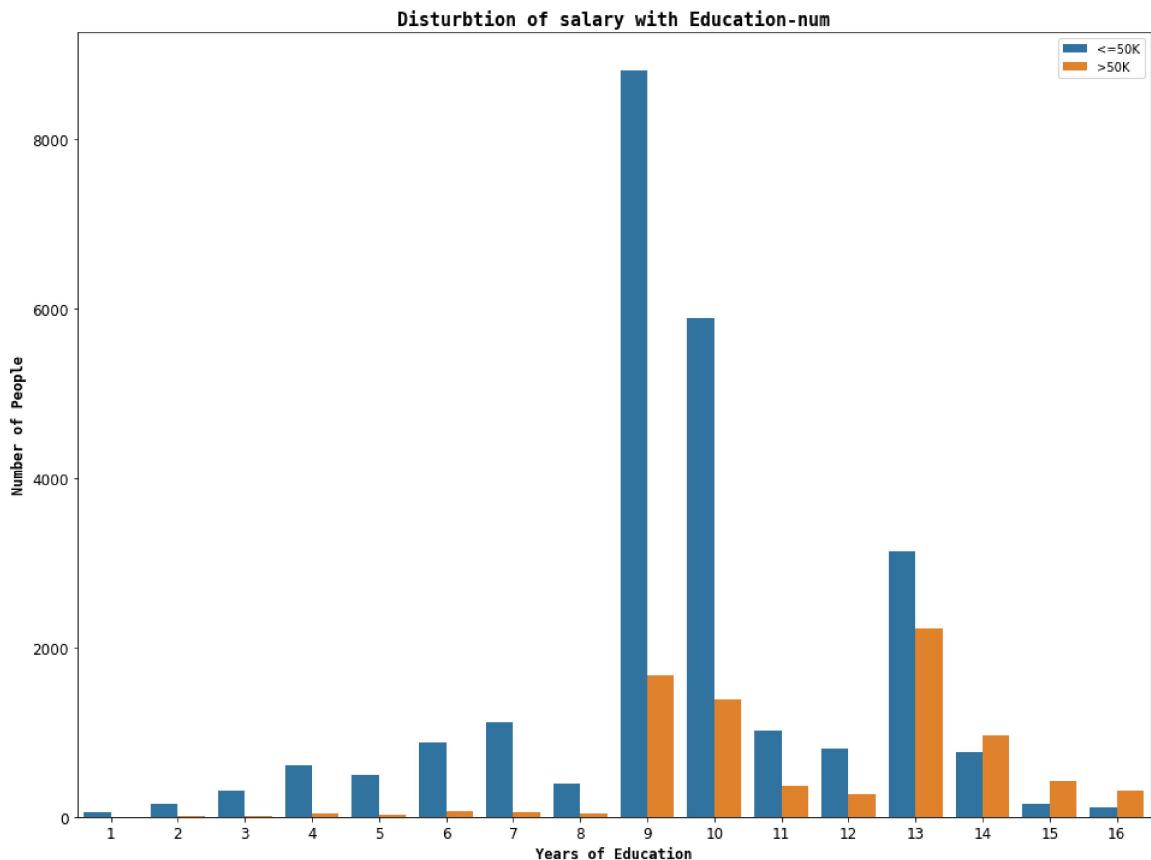
In [30]:

```
1 # countplot of salary by education
2
3 plt.figure(figsize=(20,15))
4 sns.countplot(df['education'],
5                 hue=df['salary'], palette='dark')
6 plt.title('Distruibution of Salary with Education', fontdict={
7     'fontname':'Monospace', 'fontsize':13, 'fontweight':'bold'})
8 plt.xlabel('Education', fontdict={
9     'fontname':'Monospace', 'fontsize':10, 'fontweight':'bold'})
10 plt.ylabel('Number of people',fontdict={
11     'fontname':'Monospace', 'fontsize':10, 'fontweight':'bold'})
12 plt.legend(loc=1, prop={'size':15})
13 plt.tick_params(labelsize=10)
14 plt.show()
```



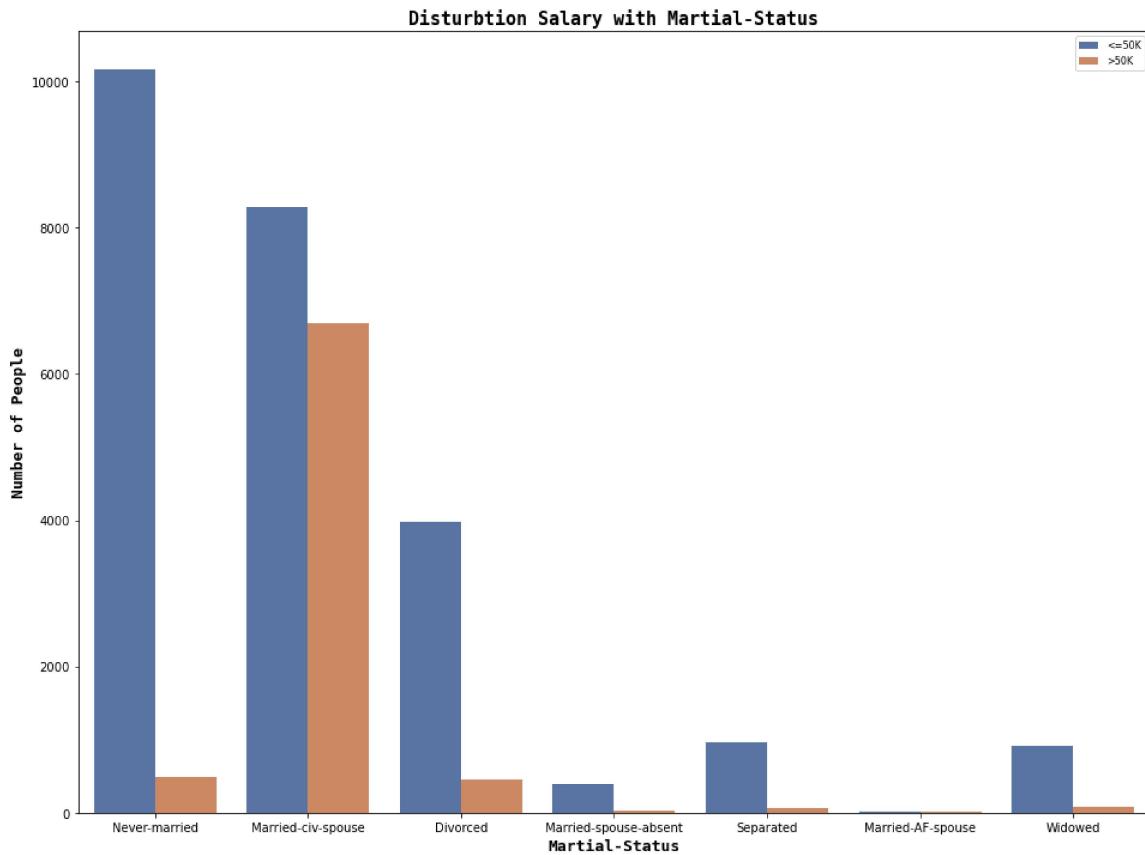
In [31]:

```
1 # Countplot of salary with years of Education-num
2
3 plt.figure(figsize=(16,12))
4 sns.countplot(df['education-num'], hue=df['salary'])
5 plt.title('Disturbtion of salary with Education-num', fontdict={
6     'fontname':'Monospace', 'fontsize':15, 'fontweight':'bold'})
7 plt.xlabel('Years of Education', fontdict={
8     'fontname':'Monospace', 'fontsize':12, 'fontweight':'bold'})
9 plt.ylabel('Number of People', fontdict={
10    'fontname':'Monospace', 'fontsize':12, 'fontweight':'bold'})
11 plt.legend(loc=1, prop={'size':10})
12 plt.tick_params(labelsize=12)
13 plt.show()
```



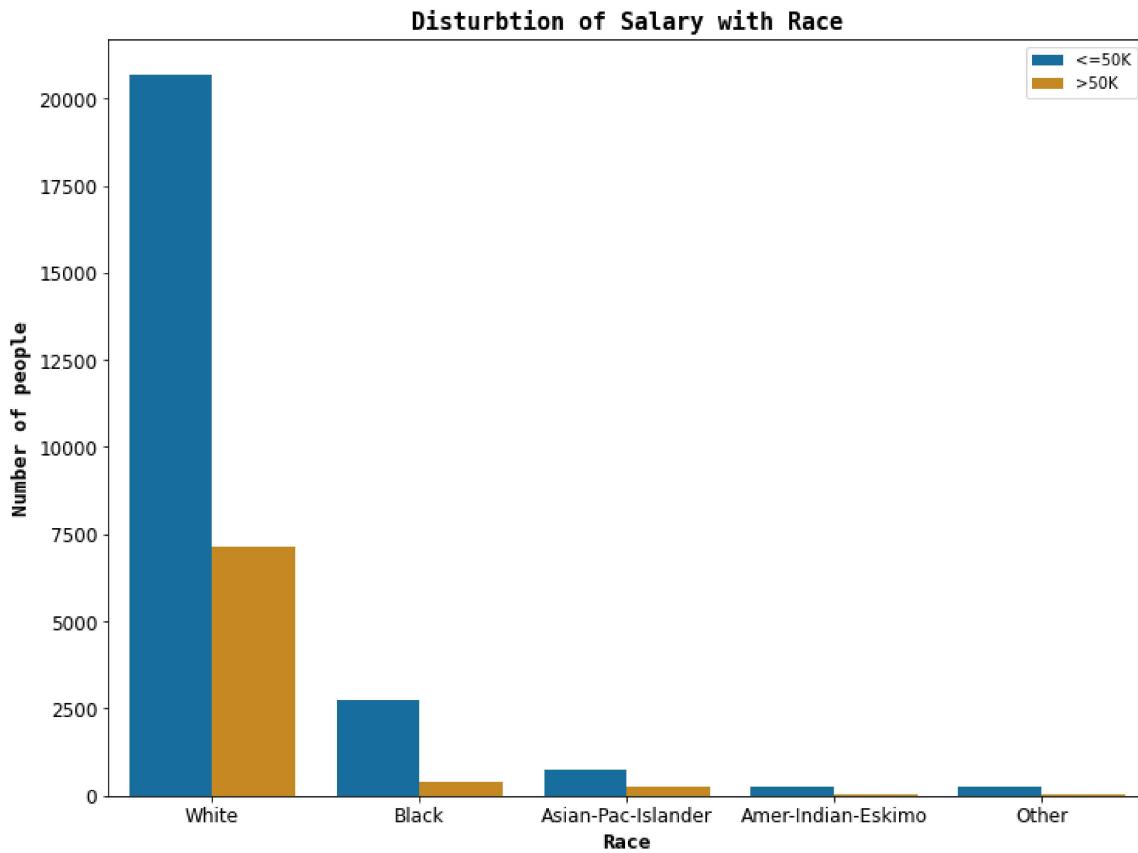
In [32]:

```
1 # creating of salary with Marital status
2
3 plt.figure(figsize=(16,12))
4 sns.countplot(df['marital-status'],
5                 hue=df['salary'], palette='deep')
6 plt.title('Disturbation Salary with Martial-Status', fontdict={
7     'fontname':'Monospace', 'fontsize':15, 'fontweight':'bold'})
8 plt.xlabel('Martial-Status', fontdict={
9     'fontname':'Monospace', 'fontsize':13, 'fontweight':'bold'})
10 plt.ylabel('Number of People', fontdict={
11     'fontname':'Monospace', 'fontsize':13, 'fontweight':'bold'})
12 plt.legend(loc=1, prop={'size':8})
13 plt.tick_params(labelsize=10)
14 plt.show()
```



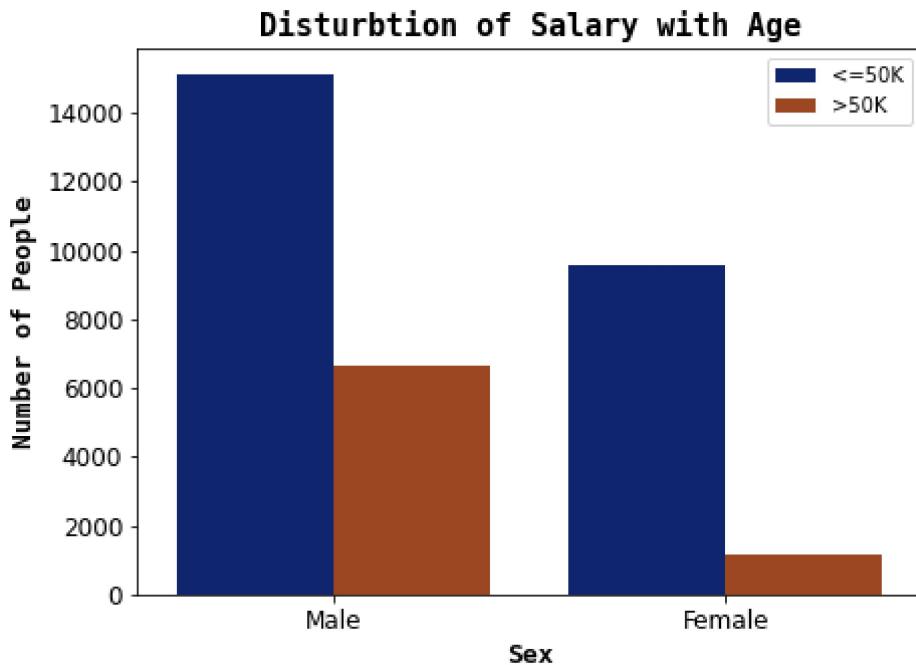
In [33]:

```
1 # creating countplot of salary with race
2
3 plt.figure(figsize=(12,9))
4 sns.countplot(df['race'],
5                 hue=df['salary'], palette='colorblind')
6 plt.title('Disturbtion of Salary with Race', fontdict={
7     'fontname':'Monospace', 'fontsize':15, 'fontweight':'bold'})
8 plt.xlabel('Race', fontdict={
9     'fontname':'Monospace', 'fontsize':13, 'fontweight':'bold'})
10 plt.ylabel('Number of people', fontdict={
11     'fontname':'Monospace', 'fontsize':13, 'fontweight':'bold'})
12 plt.legend(loc=1, prop={'size':10})
13 plt.tick_params(labelsize=12)
14 plt.show()
```



In [34]:

```
1 # creating countplot of salary with sex
2
3 plt.figure(figsize=(7,5))
4 sns.countplot(df['sex'],
5                 hue=df['salary'], palette='dark')
6 plt.title('Disturbtion of Salary with Age', fontdict={
7     'fontname':'Monospace', 'fontsize':15, 'fontweight':'bold'})
8 plt.xlabel('Sex', fontdict={
9     'fontname':'Monospace', 'fontsize':13, 'fontweight':'bold'})
10 plt.ylabel('Number of People', fontdict={
11     'fontname':'Monospace', 'fontsize':13, 'fontweight':'bold'})
12 plt.legend(loc=1, prop={'size':10})
13 plt.tick_params(labelsize=12)
14 plt.show()
```



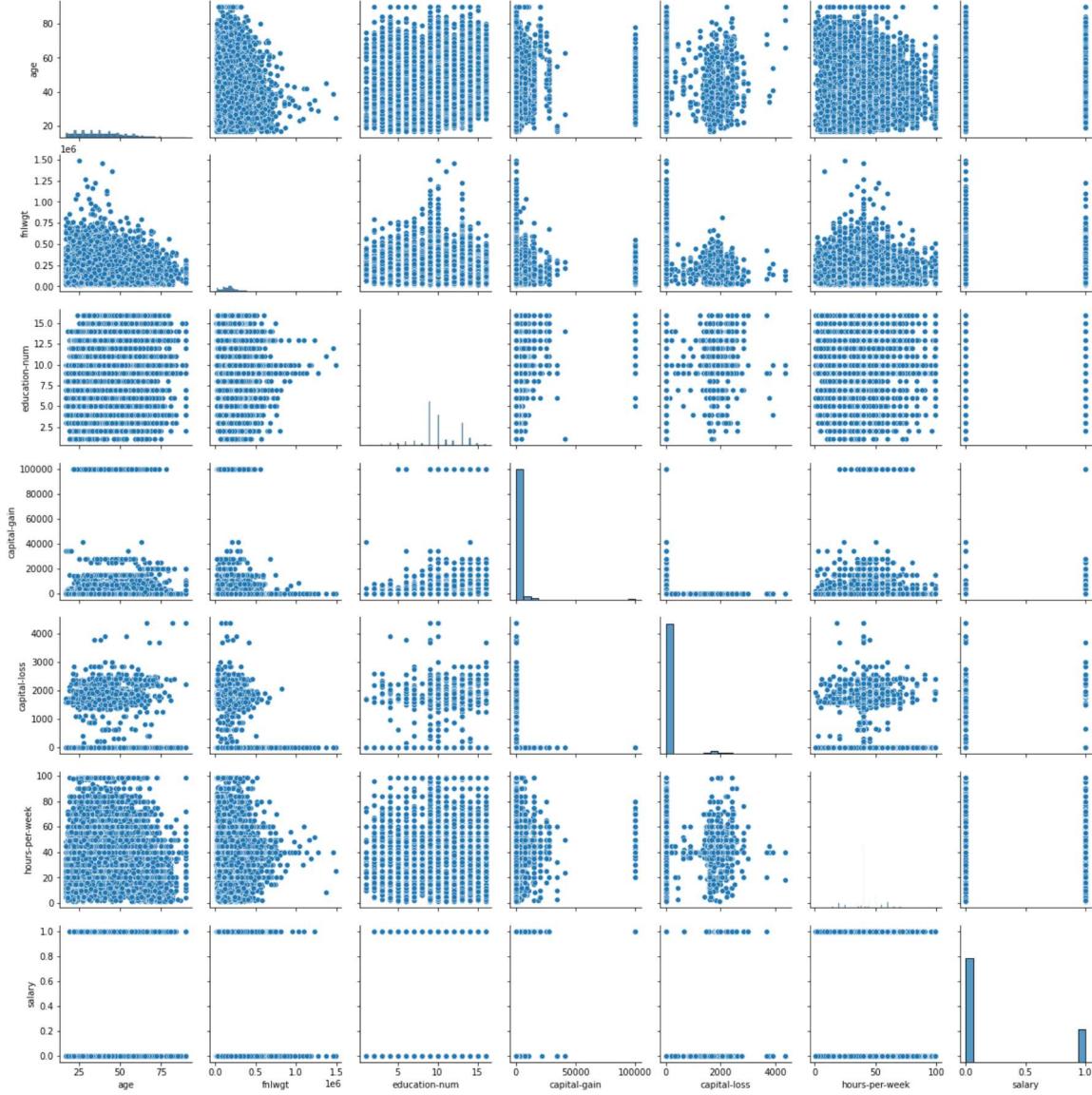
Multiple Plot with help of Pairplot

In [35]:

```
1 from sklearn.preprocessing import LabelEncoder
2 le = LabelEncoder()
3 df['salary'] = le.fit_transform(df['salary'])
```

In [36]:

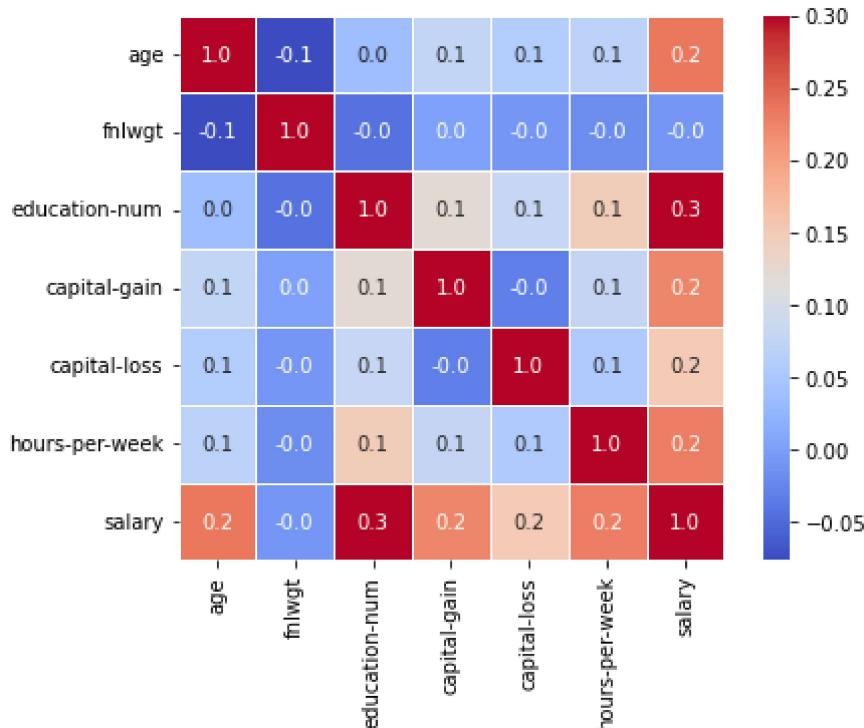
```
1 # pairplot of df (adult dataset)
2 sns.pairplot(df)
3 plt.show()
```



Correlation

In [37]:

```
1 # Correlation matrix
2
3 corr = df.corr()
4 f,ax = plt.subplots(figsize=(7,5))
5 sns.heatmap(corr, cbar = True, square = True, annot = True, fmt= '.1f',
6               xticklabels= True, yticklabels= True, vmax=0.3,
7               cmap="coolwarm", linewidths=.5, ax=ax)
8 plt.show()
```



LabelEncoder

In [38]:

```
1 for col in df.columns:
2     if df[col].dtypes == 'object':
3         encoder = LabelEncoder()
4         df[col] = encoder.fit_transform(df[col])
```

In [39]:

```
1 X = df.drop('salary',axis=1)
2 Y = df['salary']
```

```
1 for col in X.columns:
2     scaler = StandardScaler()
3     X[col] = scaler.fit_transform(X[col].values.reshape(-1,1))
```

Modelling

In [40]:

```
1 from imblearn.over_sampling import RandomOverSampler  
2 ros = RandomOverSampler(random_state=42)
```

In [41]:

```
1 ros.fit(X, Y)
```

Out[41]:

```
▶ RandomOverSampler
```

In [42]:

```
1 X_resampled, Y_resampled = ros.fit_resample(X, Y)
```

In [43]:

```
1 round(Y_resampled.value_counts(normalize=True) * 100, 2).astype('str') + '%'
```

Out[43]:

```
0    50.0 %  
1    50.0 %  
Name: salary, dtype: object
```

In [44]:

```
1 from sklearn.model_selection import train_test_split  
2 X_train, X_test, Y_train, Y_test = train_test_split(  
3     X_resampled, Y_resampled, test_size=0.2, random_state=42)
```

In [45]:

```
1 print("X_train shape:", X_train.shape)  
2 print("X_test shape:", X_test.shape)  
3 print("Y_train shape:", Y_train.shape)  
4 print("Y_test shape:", Y_test.shape)
```

```
X_train shape: (39516, 14)  
X_test shape: (9880, 14)  
Y_train shape: (39516,)  
Y_test shape: (9880,)
```

Logistic Regression

In [46]:

```
1 from sklearn.linear_model import LogisticRegression  
2 log_reg = LogisticRegression(random_state=42)
```

In [47]:

```
1 log_reg.fit(X_train, Y_train)
```

Out[47]:

```
► LogisticRegression
```

In [48]:

```
1 Y_pred_log_reg = log_reg.predict(X_test)
```

In [49]:

```
1 print("Logistic Regression:")
2 print('Accuracy Score: ',
3      round(accuracy_score(Y_test, Y_pred_log_reg)* 100, 2))
4 print('F1 Score: ',
5      round(f1_score(Y_test, Y_pred_log_reg) * 100, 2))
```

Logistic Regression:

Accuracy Score: 65.37

F1 Score: 63.23

KNN Classifier

In [50]:

```
1 from sklearn.neighbors import KNeighborsClassifier
2 knn = KNeighborsClassifier()
```

In [51]:

```
1 knn.fit(X_train, Y_train)
```

Out[51]:

```
► KNeighborsClassifier
```

In [52]:

```
1 Y_pred_knn = knn.predict(X_test)
```

In [53]:

```
1 print('KNN Classifier: ')
2 print('Accuracy Score: ',
3      round(accuracy_score(Y_test, Y_pred_knn) * 100, 2))
4 print('F1 Score: ',
5      round(f1_score(Y_test, Y_pred_knn) * 100, 2))
```

KNN Classifier:

Accuracy Score: 72.75

F1 Score: 74.82

Support Vector Classifier (SVC)

In [54]:

```
1 from sklearn.svm import SVC  
2 svc = SVC(random_state=42)
```

In [55]:

```
1 svc.fit(X_train, Y_train)
```

Out[55]:

```
▶ SVC
```

In [56]:

```
1 Y_pred_svc = svc.predict(X_test)
```

In [57]:

```
1 print('Support Vector Classifier: ')  
2 print('Accuracy Score: ',  
3       round(accuracy_score(Y_test, Y_pred_svc) * 100 , 2))  
4 print('F1 Score: ',  
5       round(f1_score(Y_test, Y_pred_svc) * 100 , 2))
```

Support Vector Classifier:

Accuracy Score: 59.23

F1 Score: 32.01

Naive Bayes Classifier

In [58]:

```
1 from sklearn.naive_bayes import GaussianNB  
2 nb = GaussianNB()
```

In [59]:

```
1 nb.fit(X_train, Y_train)
```

Out[59]:

```
▶ GaussianNB
```

In [60]:

```
1 Y_pred_nb = nb.predict(X_test)
```

In [61]:

```
1 print('Support Vector Classifier: ')
2 print('Accuracy Score: ',
3      round(accuracy_score(Y_test, Y_pred_nb) * 100, 2))
4 print('F1 Score: ',
5      round(f1_score(Y_test, Y_pred_nb) * 100 , 2))
```

Support Vector Classifier:

Accuracy Score: 62.59

F1 Score: 44.98

Decision Tree Classifier

In [62]:

```
1 from sklearn.tree import DecisionTreeClassifier
2 dec_tree = DecisionTreeClassifier(random_state=42)
```

In [63]:

```
1 dec_tree.fit(X_train, Y_train)
```

Out[63]:

```
▶ DecisionTreeClassifier
```

In [64]:

```
1 Y_pred_dec_tree = dec_tree.predict(X_test)
```

In [65]:

```
1 print('Support Vector Classifier: ')
2 print('Accuracy Score: ',
3      round(accuracy_score(Y_test, Y_pred_dec_tree)* 100, 2))
4 print('F1 Score: ',
5      round(f1_score(Y_test, Y_pred_dec_tree) * 100, 2))
```

Support Vector Classifier:

Accuracy Score: 91.13

F1 Score: 91.57

Random Forest Classifier

In [66]:

```
1 from sklearn.ensemble import RandomForestClassifier
2 rf = RandomForestClassifier(random_state=42)
```

In [67]:

```
1 rf.fit(X_train, Y_train)
```

Out[67]:

```
► RandomForestClassifier
```

In [68]:

```
1 Y_pred_rf = rf.predict(X_test)
```

In [69]:

```
1 print('Support Vector Classifier: ')
2 print('Accuracy Score: ',
3      round(accuracy_score(Y_test, Y_pred_rf) * 100 , 2))
4 print('F1 Score: ',
5      round(f1_score(Y_test, Y_pred_rf)* 100 , 2))
```

Support Vector Classifier:

Accuracy Score: 92.98

F1 Score: 93.23

XGB Classifier

In [70]:

```
1 from xgboost import XGBClassifier
2 xgb = XGBClassifier()
```

In [71]:

```
1 xgb.fit(X_train, Y_train)
```

Out[71]:

```
► XGBClassifier
```

In [72]:

```
1 Y_pred_xgb = xgb.predict(X_test)
```

In [73]:

```
1 print('Support Vector Classifier: ')
2 print('Accuracy Score: ',
3      round(accuracy_score(Y_test, Y_pred_xgb) * 100 , 2))
4 print('F1 Score: ',
5      round(f1_score(Y_test, Y_pred_xgb) * 100 , 2))
```

Support Vector Classifier:

Accuracy Score: 86.95

F1 Score: 87.33

AdaBoost Classifier

In [74]:

```
1 from sklearn.ensemble import AdaBoostClassifier  
2 adaboost = AdaBoostClassifier(random_state=42)
```

In [75]:

```
1 adaboost.fit(X_train, Y_train)
```

Out[75]:

```
▶ AdaBoostClassifier
```

In [76]:

```
1 Y_pred_adaboost = adaboost.predict(X_test)
```

In [77]:

```
1 print('Support Vector Classifier: ')  
2 print('Accuracy Score: ',  
3       round(accuracy_score(Y_test, Y_pred_adaboost) * 100, 2))  
4 print('F1 Score: ',  
5       round(f1_score(Y_test, Y_pred_adaboost) * 100, 2))
```

Support Vector Classifier:

Accuracy Score: 82.18

F1 Score: 82.42

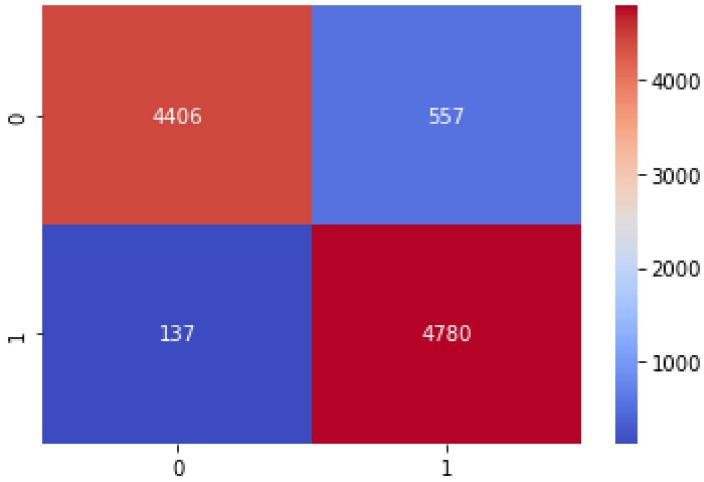
Confusion Matrix

In [78]:

```
1 from sklearn.metrics import confusion_matrix  
2 con_mat = confusion_matrix(Y_test, Y_pred_rf)
```

In [79]:

```
1 sns.heatmap(con_mat, annot=True, cmap='coolwarm', fmt='d')
2 plt.show()
```



Classification Report

In [80]:

```
1 from sklearn.metrics import classification_report
2 print(classification_report(Y_test, Y_pred_rf))
```

	precision	recall	f1-score	support
0	0.97	0.89	0.93	4963
1	0.90	0.97	0.93	4917
accuracy			0.93	9880
macro avg	0.93	0.93	0.93	9880
weighted avg	0.93	0.93	0.93	9880

In []:

```
1
```