

## Import libraries

```
1. import pandas as pd
2. from scipy import stats
3. import numpy as np
4. import seaborn as sns
5. import matplotlib.pyplot as plt
6. from imblearn.combine import SMOTEENN
7. from sklearn.linear_model import LogisticRegression
8. from sklearn.tree import DecisionTreeClassifier
9. from sklearn.ensemble import RandomForestClassifier
10. from xgboost import XGBClassifier
11. from sklearn.metrics import confusion_matrix, classification_report,
    accuracy_score, roc_auc_score, roc_curve
12. from sklearn.model_selection import train_test_split, GridSearchCV, cross_val_score
13. import pickle
14. import warnings
15. warnings.filterwarnings('ignore')
16. import matplotlib
17. matplotlib.rcParams["figure.dpi"] = 80
```

## Loading the data

```
1. data_train = pd.read_csv("adult.csv", header=None,
    names=["age", "workclass", "final_weight", "education", "education-num", "marital-
    status", "occupation", "relationship", "race", "sex", "capital-gain", "capital-
    loss", "hours-per-week", "native-country", "income"])
2. data_test = pd.read_csv("adult.test", header=None,
    names=["age", "workclass", "final_weight", "education", "education-num", "marital-
    status", "occupation", "relationship", "race", "sex", "capital-gain", "capital-
    loss", "hours-per-week", "native-country", "income"])
```

## EDA

### What is the dimension of data?

```
1. def row_col(data):
2.     print("Total rows present in dataset - {}".format(data.shape[0]))
3.     print("Total columns present in dataset - {}".format(data.shape[1]))
4.
5. print("Training Data:")
6. row_col(data_train)
7. print()
8. print("Testing Data:")
9. row_col(data_test)
```

```
Training Data:
Total rows present in dataset - 32561
Total columns present in dataset - 15
```

```
Testing Data:
Total rows present in dataset - 16281
Total columns present in dataset - 15
```

## Categorical features and numerical features

```
1. def get_cat_num(data):
2.
3.     categorical = []
4.     numerical = []
5.
6.     for column in data.drop("income", axis=1).columns:
7.         if data[column].dtype == "O":
8.             categorical.append(column)
9.         else:
10.            numerical.append(column)
11.     return categorical, numerical
12.
13. categorical, numerical = get_cat_num(data_train)
14. print("There are {} categorical features.".format(len(categorical)))
15. print("Categorical features: {}".format(categorical))
16. print()
17. print("There are {} numerical features.".format(len(numerical)))
18. print("Numerical features: {}".format(numerical))
```

There are 8 categorical features.

Categorical features: ['workclass', 'education', 'marital-status', 'occupation', 'relationship', 'race', 'sex', 'native-country']

There are 6 numerical features.

Numerical features: ['age', 'final\_weight', 'education-num', 'capital-gain', 'capital-loss', 'hours-per-week']

```
1. def strip_data(data):
2.     for cat_f in categorical:
3.         data[cat_f] = data[cat_f].str.strip()
4.
5. strip_data(data_train)
6. strip_data(data_test)
```

## Is there any duplicated rows?

```
1. if data_train.duplicated().sum():
2.     print("Yes! There are", data_train.duplicated().sum(), "duplicate rows.")
3.     data_train.drop_duplicates(inplace=True)
```

Yes! There are 24 duplicate rows.

## Replace '?' with 'Unknown'

```
1. data_train = data_train.replace('?', 'Unknown')
2. data_test = data_test.replace('?', 'Unknown')
```

## Analysis of Categorical Data

```
1. summary = pd.DataFrame(columns=["Feature", "Relationship", "Strength"])
2.
3. def get_plot(feature, rot=0):
4.     crossTable = pd.crosstab(data_train[feature], data_train["income"])
```

```

5.     percTable = crossTable.div(crossTable.sum(axis=0), axis=1)*100
6.     ax = percTable.plot(kind="bar", figsize=(15,5), title="Income of Adult -
    "+feature, rot=rot)
7.     ax.set_ylabel("Percentage", fontsize=10)
8.     ax.set_xlabel(feature, fontsize=12)
9.     plt.show()
10.
11. def cramers_v_test(feature):
12.     global summary
13.     crossTable = pd.crosstab(data_train[feature], data_train["income"])
14.     chiVal, pVal, df, exp = stats.chi2_contingency(crossTable)
15.     print("Chi-square value :", round(chiVal,4))
16.     print("P Value          :", pVal)
17.     print()
18.
19.     if pVal < 0.05:
20.         r = len(crossTable.index)
21.         c = len(crossTable.columns)
22.         n = crossTable.to_numpy().sum()
23.         v = np.sqrt(chiVal / (n * (min(r,c)-1)))
24.
25.         print("Income is dependent on", feature, end=". ")
26.         print("Cramér's V: ", round(v,2))
27.
28.         if df == 1:
29.             if v < 0.10:
30.                 relationship = "Negligible"
31.                 print("Relationship: negligible")
32.             elif v < 0.30:
33.                 print("Relationship: small")
34.             elif v < 0.50:
35.                 print("Relationship: medium")
36.             else:
37.                 print("Relationship: large")
38.
39.         elif df == 2:
40.             if v < 0.07:
41.                 relationship = "Negligible"
42.                 print("Relationship: negligible")
43.             elif v < 0.21:
44.                 relationship = "Small"
45.                 print("Relationship: small")
46.             elif v < 0.35:
47.                 relationship = "Medium"
48.                 print("Relationship: medium")
49.             else:
50.                 relationship = "Large"
51.                 print("Relationship: large")
52.
53.         elif df == 3:
54.             if v < 0.06:
55.                 relationship = "Negligible"
56.                 print("Relationship: negligible")
57.             elif v < 0.17:
58.                 relationship = "Small"
59.                 print("Relationship: small")
60.             elif v < 0.29:
61.                 relationship = "Medium"
62.                 print("Relationship: medium")
63.             else:
64.                 relationship = "Large"
65.                 print("Relationship: large")
66.
67.         elif df == 4:
68.             if v < 0.05:
69.                 relationship = "Negligible"

```

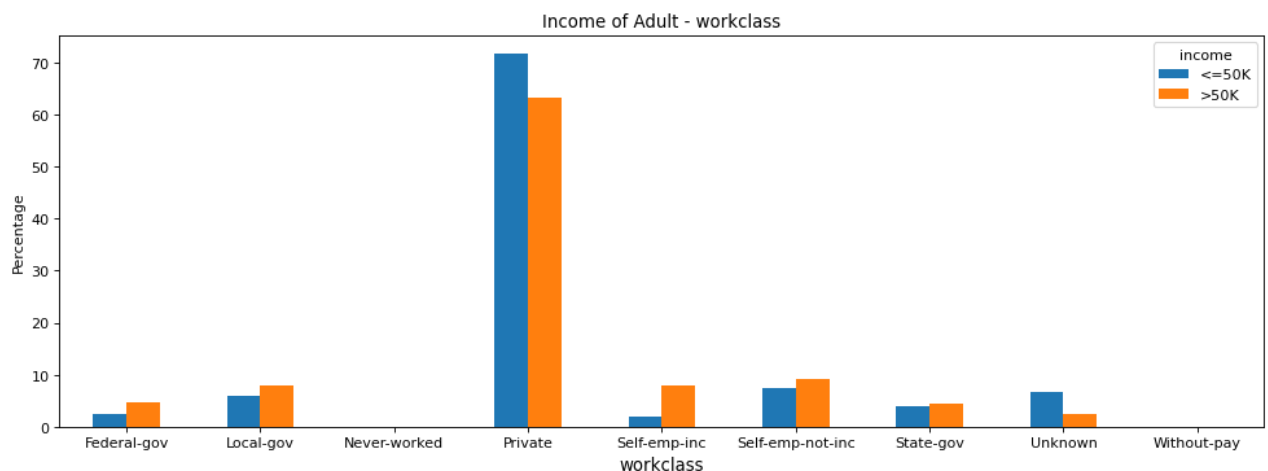
```

70.         print("Relationship: negligible")
71.     elif v < 0.15:
72.         relationship = "Small"
73.         print("Relationship: small")
74.     elif v < 0.25:
75.         relationship = "Medium"
76.         print("Relationship: medium")
77.     else:
78.         relationship = "Large"
79.         print("Relationship: large")
80.
81.     else:
82.         if v < 0.05:
83.             relationship = "Negligible"
84.             print("Relationship: negligible")
85.         elif v < 0.13:
86.             relationship = "Small"
87.             print("Relationship: small")
88.         elif v < 0.22:
89.             relationship = "Medium"
90.             print("Relationship: medium")
91.         else:
92.             relationship = "Large"
93.             print("Relationship: large")
94.
95.     summary = summary.append({"Feature":feature, "Relationship":relationship,
96. "Strength":round(v,2)}, ignore_index=True)
97.     else:
98.         print("Income is not dependent on",feature)

```

## Workclass

```
1. get_plot("workclass")
```



```
1. cramers_v_test("workclass")
```

```

Chi-square value : 1044.6962
P Value          : 3.352256069028484e-220

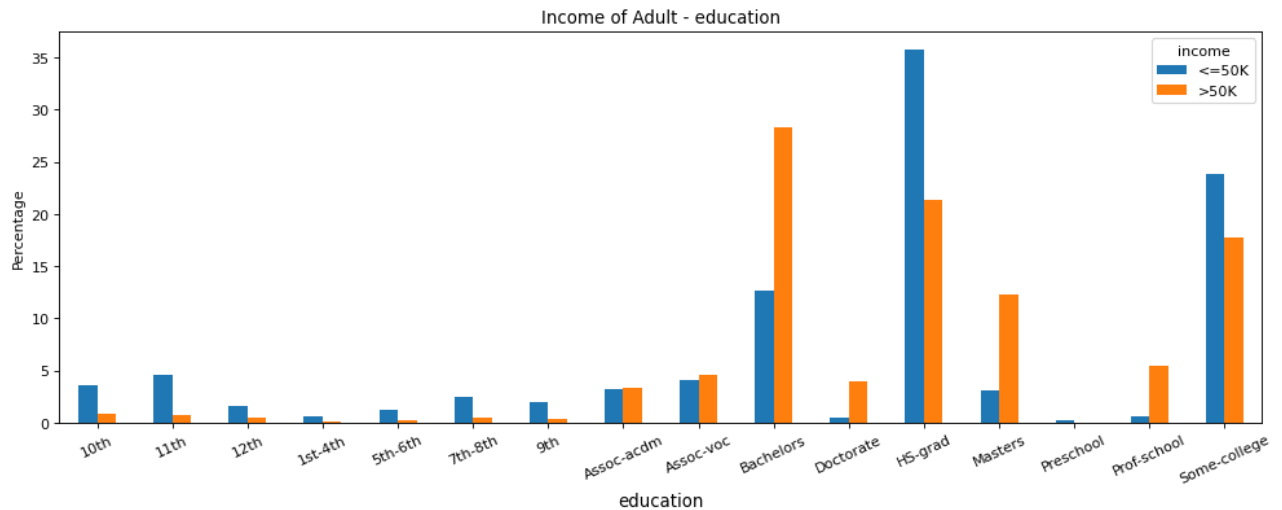
```

Income is dependent on workclass. Cramér's V: 0.18

Relationship: medium

## Education

```
1. get_plot("education", 25)
```



```
1. data_train["education"] = np.where(data_train["education"].isin(["Preschool", "1st-4th", "5th-6th", "7th-8th", "9th", "10th", "11th", "12th"]), "School", data_train["education"])
2. data_train['education'] = data_train['education'].map({'School':0, 'HS-grad':1, 'Some-college':2, 'Assoc-voc':3, 'Assoc-acdm':4, 'Bachelors':5, 'Masters':6, 'Prof-school':7, 'Doctorate':8})
3. data_test["education"] = np.where(data_test["education"].isin(["Preschool", "1st-4th", "5th-6th", "7th-8th", "9th", "10th", "11th", "12th"]), "School", data_test["education"])
4. data_test['education'] = data_test['education'].map({'School':0, 'HS-grad':1, 'Some-college':2, 'Assoc-voc':3, 'Assoc-acdm':4, 'Bachelors':5, 'Masters':6, 'Prof-school':7, 'Doctorate':8})
```

```
1. cramers_v_test("education")
```

Chi-square value : 4425.2842

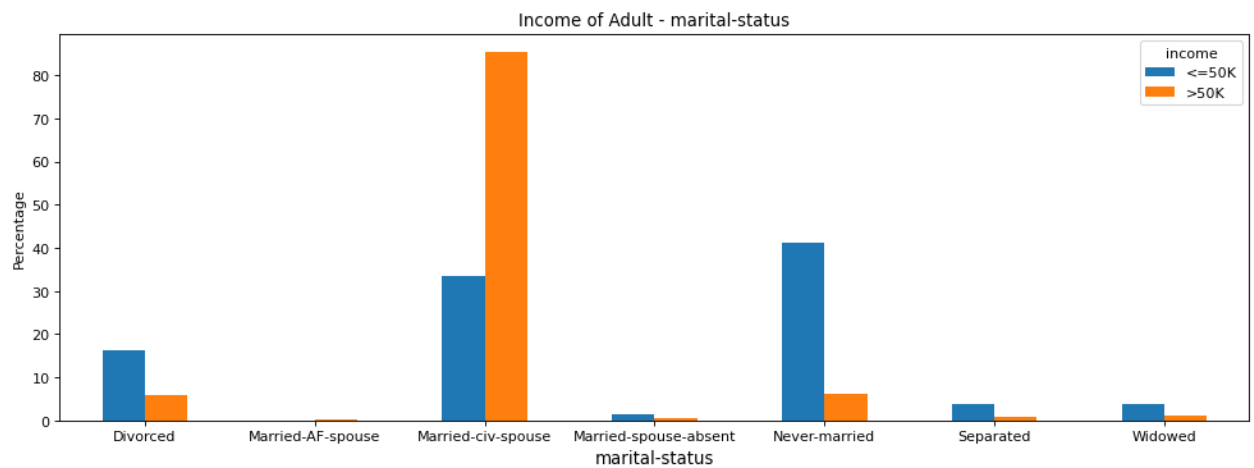
P Value : 0.0

Income is dependent on education. Cramér's V: 0.37

Relationship: large

## Marital-status

```
1. get_plot("marital-status")
```



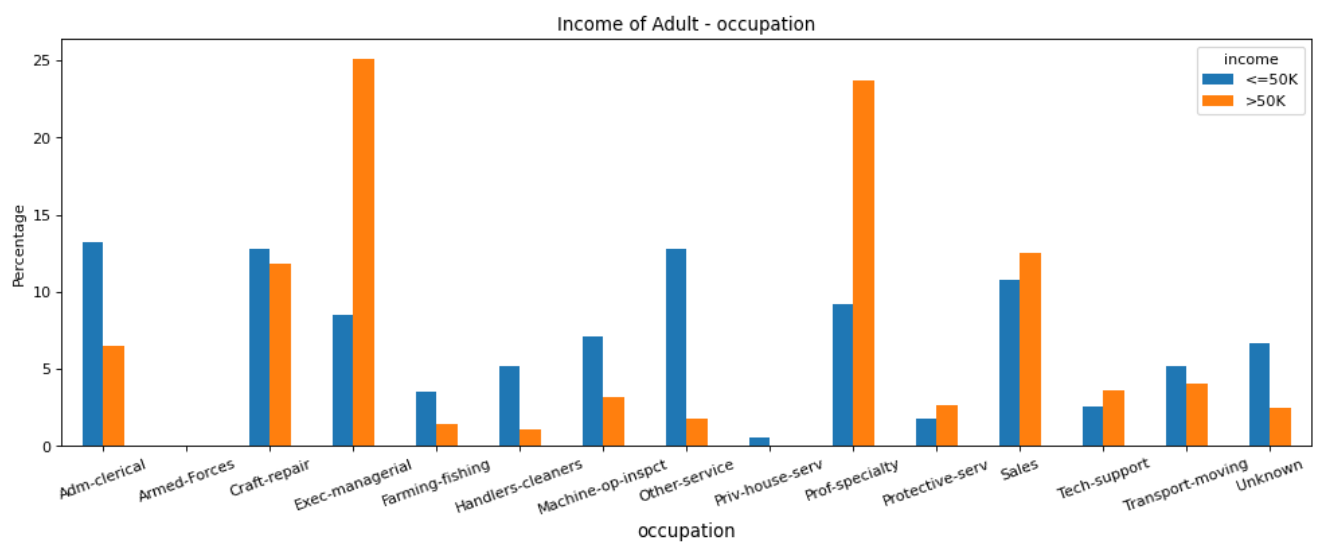
```
1. cramers_v_test("marital-status")
```

Chi-square value : 6510.3321  
P Value : 0.0

Income is dependent on marital-status. Cramér's V: 0.45  
Relationship: large

## Occupation

```
1. get_plot("occupation", 20)
```



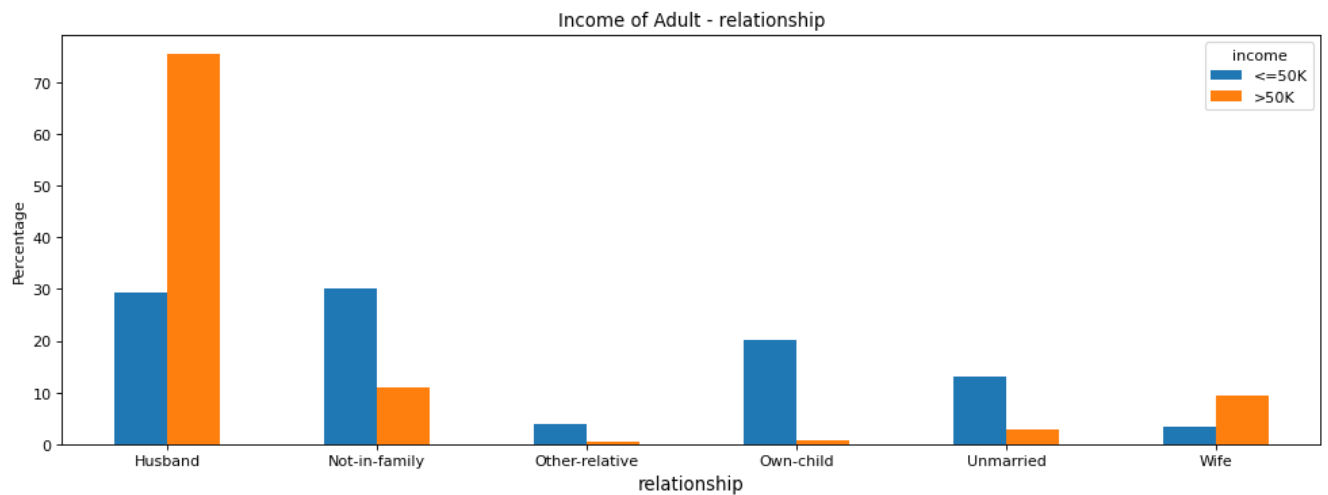
```
1. cramers_v_test("occupation")
```

Chi-square value : 4030.2092  
P Value : 0.0

Income is dependent on occupation. Cramér's V: 0.35  
Relationship: large

## Relationship

```
1. get_plot("relationship")
```



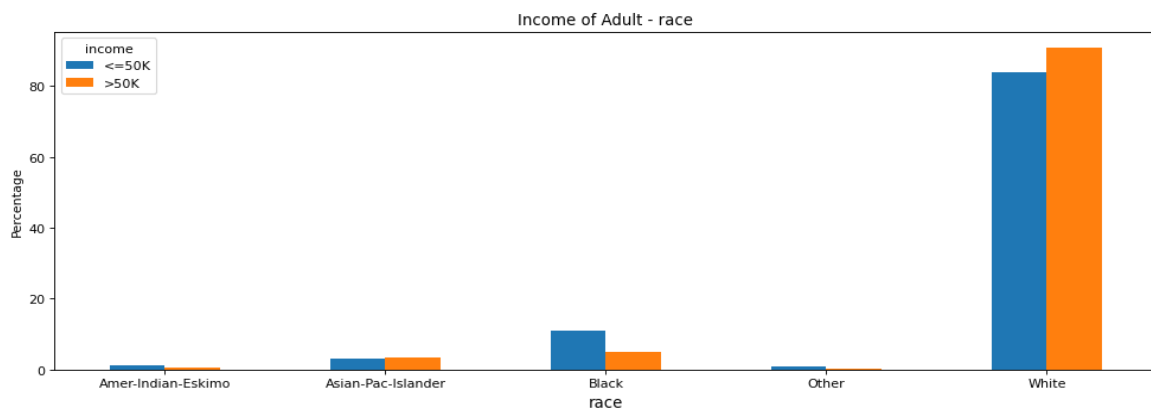
```
1. cramers_v_test("relationship")
```

Chi-square value : 6692.0988  
P Value : 0.0

Income is dependent on relationship. Cramér's V: 0.45  
Relationship: large

## Race

```
1. get_plot("race")
```



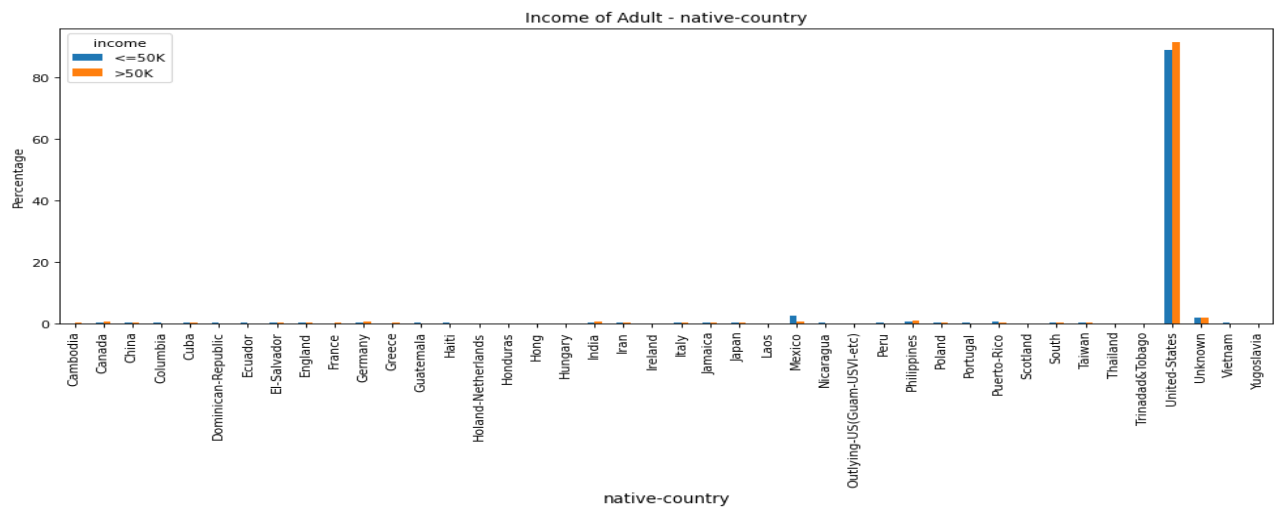
```
1. cramers_v_test("race")
```

```
Chi-square value : 330.9434
P Value          : 2.2797874171824478e-70
```

```
Income is dependent on race. Cramér's V: 0.1
Relationship: small
```

## Native-country

```
1. get_plot("native-country", 90)
```



```
1. cramers_v_test("native-country")
```

```
Chi-square value : 315.4485
P Value          : 4.833085519399296e-44
```

```
Income is dependent on native-country. Cramér's V: 0.1
Relationship: small
```

```
1. summary
```

	Feature	Relationship	Strength
0	workclass	Medium	0.18
1	education	Large	0.37
2	marital-status	Large	0.45
3	occupation	Large	0.35
4	relationship	Large	0.45
5	race	Small	0.10
6	native-country	Small	0.10

Okay! You can clearly see **race** and **native-country** both can be dropped as they have small relationship with dependent variable.

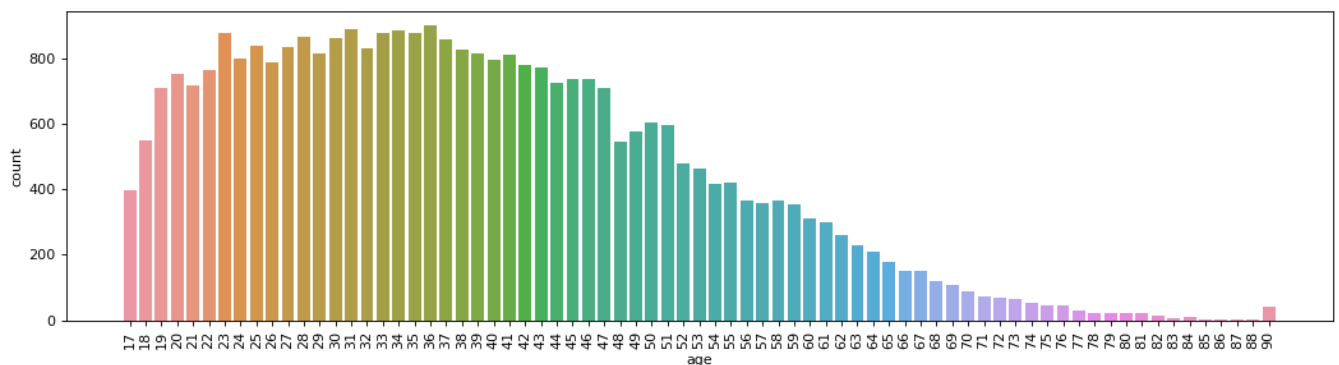


```
1. data_train.drop(["race", "native-country"], axis=1, inplace=True)
2. data_test.drop(["race", "native-country"], axis=1, inplace=True)
```

## Analysis of Numerical Data

### Age

```
1. plt.figure(figsize=(16,4))
2. sns.countplot(data_train["age"])
3. plt.xticks(rotation=90)
4. plt.plot()
```



### Education-num

```
1. data_train['education-num'].unique()
```

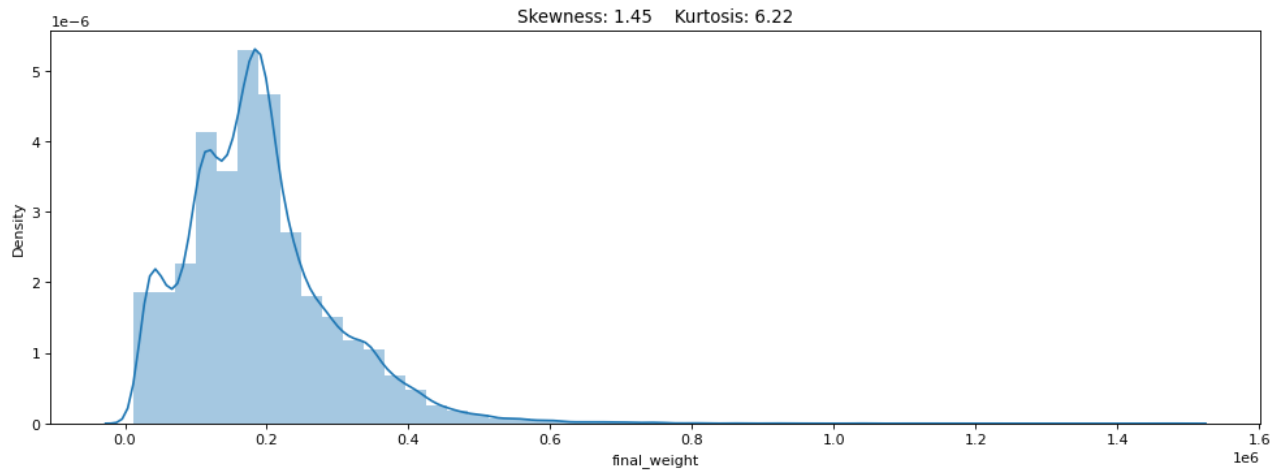
```
array([13,  9,  7, 14,  5, 10, 12, 11,  4, 16, 15,  3,  6,  2,  1,  8],
      dtype=int64)
```

Seems like **education** and **education-num** are same features. education-num is just numerical representation of education. So, **We'll drop education-num**.

```
1. data_train.drop("education-num", axis=1, inplace=True)
2. data_test.drop("education-num", axis=1, inplace=True)
```

### Final Weight

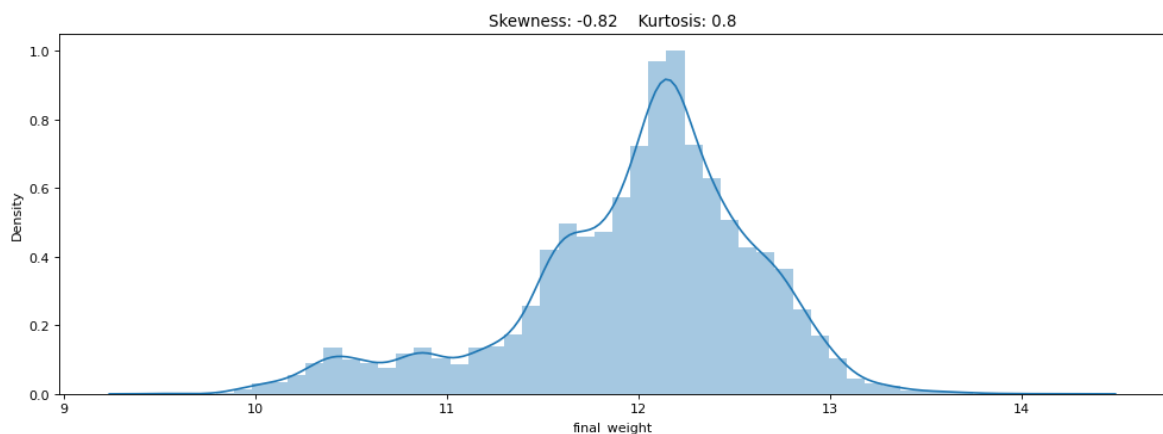
```
1. plt.figure(figsize=(15,5))
2. sns.distplot(data_train['final_weight'])
3. plt.title("Skewness: " + str(round(data_train['final_weight'].skew(),2)) + " " +
  "Kurtosis: " + str(round(data_train['final_weight'].kurtosis(),2)))
4. plt.show()
```



**Final weight** is right skewed. We need to make it normal for further tasks.

```
1. data_train['final_weight'] = np.log1p(data_train[['final_weight']])
2. data_test['final_weight'] = np.log1p(data_test[['final_weight']])
```

```
1. plt.figure(figsize=(15,5))
2. sns.distplot(data_test['final_weight'])
3. plt.title("Skewness: " + str(round(data_test['final_weight'].skew(),2)) + "      " +
            "Kurtosis: " + str(round(data_test['final_weight'].kurtosis(),2)))
4. plt.show()
```



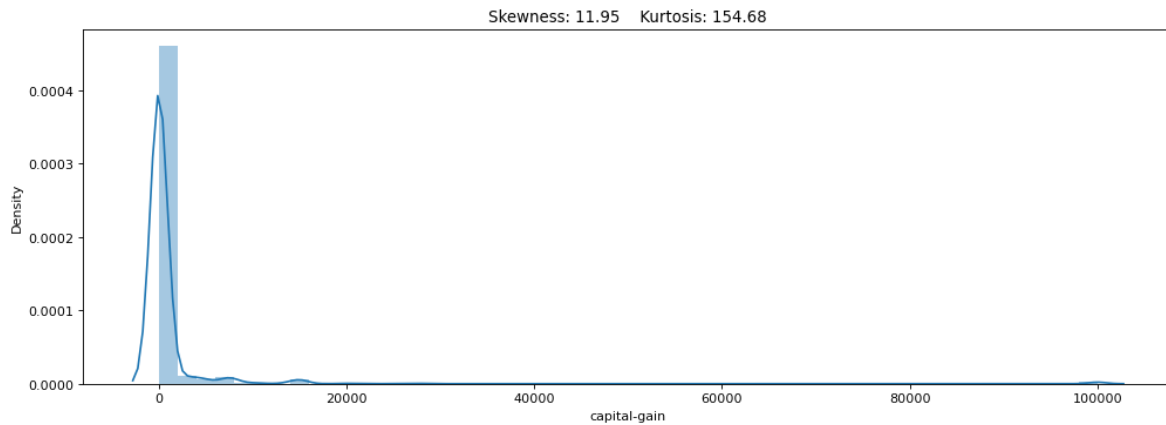
**Capital Gain**

```
1. plt.figure(figsize=(15,5))
2. sns.distplot(data_train['capital-gain'])
```

```

3. plt.title("Skewness: " + str(round(data_train['capital-gain'].skew(),2)) + " " +
    "Kurtosis: " + str(round(data_train['capital-gain'].kurtosis(),2)))
4. plt.show()

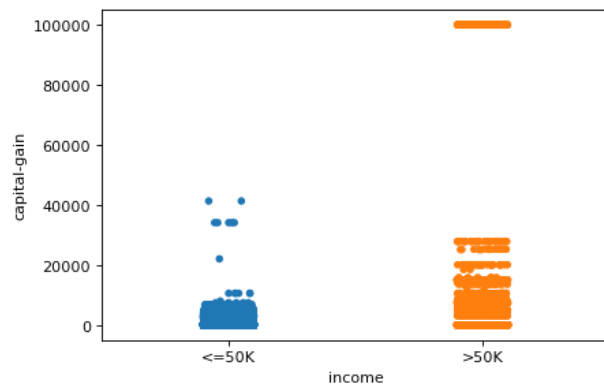
```



```

1. sns.stripplot(data=data_train, x='income', y='capital-gain')

```



```

1. data_train["capital-gain"] = np.where(data_train["capital-gain"]>=7000, 1, 0)
2. data_test["capital-gain"] = np.where(data_test["capital-gain"]>=7000, 1, 0)

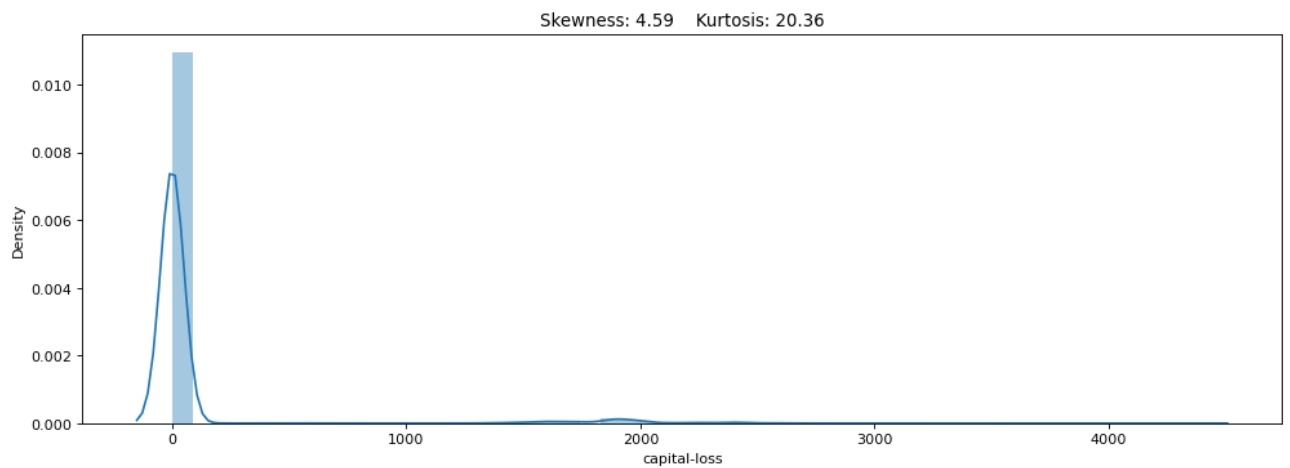
```

## Capital Loss

```

1. plt.figure(figsize=(15,5))
2. sns.distplot(data_train['capital-loss'])
3. plt.title("Skewness: " + str(round(data_train['capital-loss'].skew(),2)) + " " +
    "Kurtosis: " + str(round(data_train['capital-loss'].kurtosis(),2)))
4. plt.show()

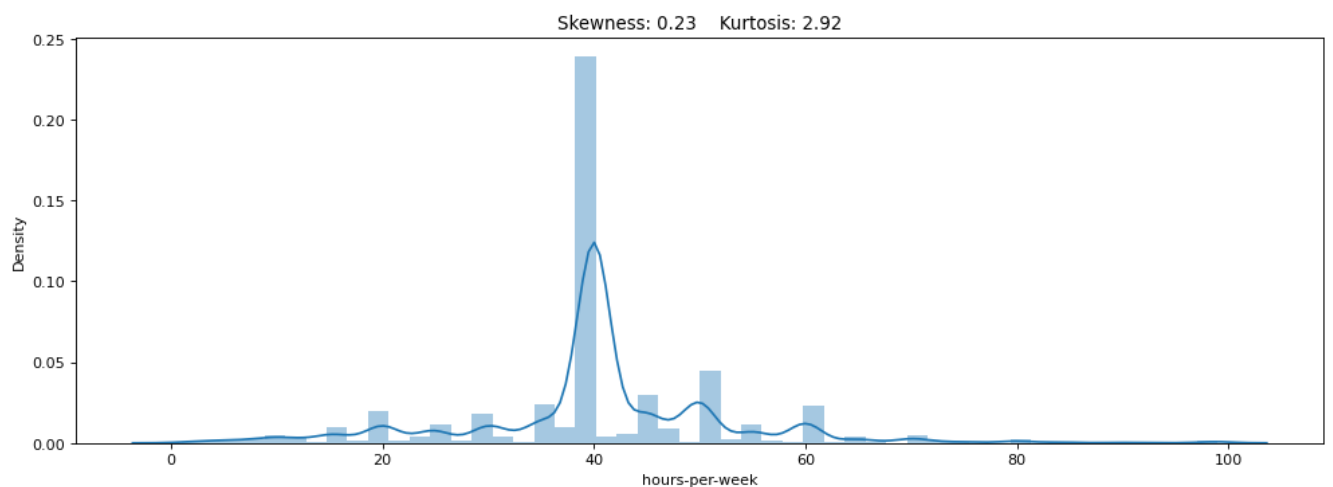
```



```
1. data_train["capital-loss"] = data_train["capital-loss"].apply(lambda x: x ** (1/2))
2. data_test["capital-loss"] = data_test["capital-loss"].apply(lambda x: x ** (1/2))
```

## Hours-per-week

```
1. plt.figure(figsize=(15,5))
2. sns.distplot(data_train['hours-per-week'])
3. plt.title("Skewness: " + str(round(data_train['hours-per-week'].skew(),2)) + "      "
            + "Kurtosis: " + str(round(data_train['hours-per-week'].kurtosis(),2)))
4. plt.show()
```

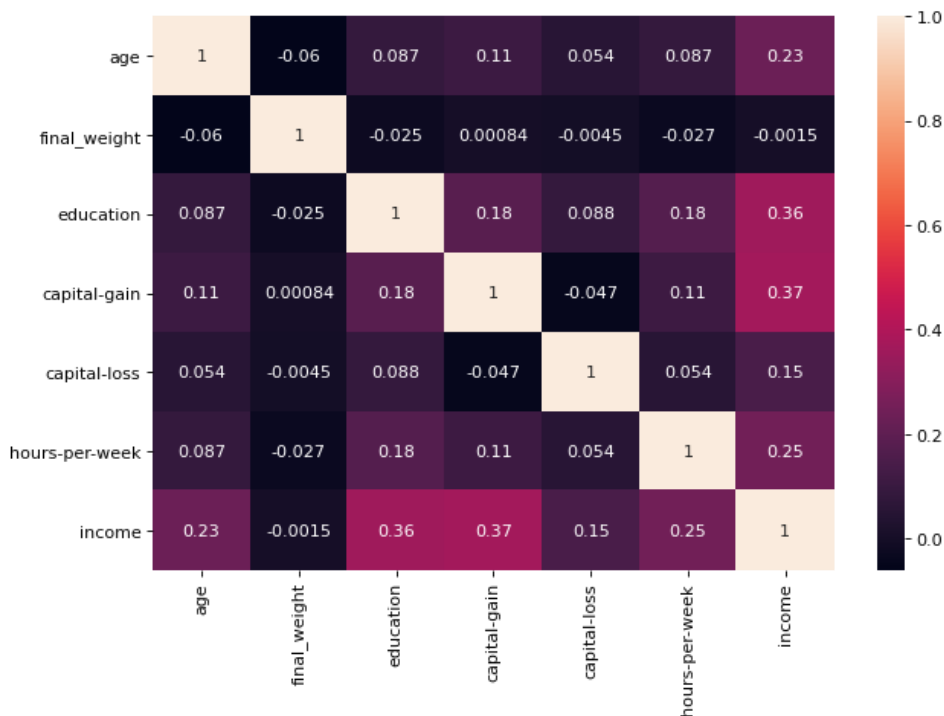


```
1. data_train['hours-per-week'] = np.where(data_train['hours-per-week']<35, 0,
np.where(data_train['hours-per-week']<50, 1, np.where(data_train['hours-per-
week']<72, 2, 3)))
2. data_test['hours-per-week'] = np.where(data_test['hours-per-week']<35, 0,
np.where(data_test['hours-per-week']<50, 1, np.where(data_test['hours-per-
week']<72, 2, 3)))
```

```
1. data_train['income'] = data_train['income'].map({"<=50K":0, ">50K":1})
2. data_test['income'] = data_test['income'].map({"<=50K.":0, ">50K.":1})
```

```
1. plt.figure(figsize=(9,6))
```

```
2. sns.heatmap(data_train.corr(), annot=True)
3. plt.plot()
```



Here, we can see that **final\_weight** is not correlated with dependent variable. So, we can drop it.

```
1. data_train.drop("final_weight", axis=1, inplace=True)
2. data_test.drop("final_weight", axis=1, inplace=True)
```

```
1. data_train.head()
```

	age	workclass	education	education-num	marital-status	occupation	relationship	sex	capital-gain	capital-loss	hours-per-week	income
0	39	State-gov	5	13	Never-married	Adm-clerical	Not-in-family	Male	0	0.0	1	0
1	50	Self-emp-not-inc	5	13	Married-civ-spouse	Exec-managerial	Husband	Male	0	0.0	0	0
2	38	Private	1	9	Divorced	Handlers-cleaners	Not-in-family	Male	0	0.0	1	0
3	53	Private	0	7	Married-civ-spouse	Handlers-cleaners	Husband	Male	0	0.0	1	0
4	28	Private	5	13	Married-civ-spouse	Prof-specialty	Wife	Female	0	0.0	1	0

```
1. X_train = data_train.drop('income', axis=1)
2. y_train = data_train['income']
3.
4. X_test = data_test.drop('income', axis=1)
5. y_test = data_test['income']
```

```
1. X_train = pd.get_dummies(X_train, drop_first=True)
2. X_test = pd.get_dummies(X_test, drop_first=True)
```

## Handling Imbalanced Data

```

1. sme = SMOTEENN()
2. X_train, y_train = sme.fit_resample(X_train, y_train)

```

```

1. y_train.value_counts()

```

```

1      18088
0       16811
Name: income, dtype: int64

```

## Selection of Model

```

1. model_params = {
2.
3.     'decision_tree' : {
4.         'model' : DecisionTreeClassifier(),
5.         'params' : {
6.             'criterion': ['gini', 'entropy'],
7.             'max_depth': [18, 20, 22],
8.             'min_samples_leaf': [2, 5, 8]
9.         }
10.    },
11.
12.    'random_forest' : {
13.        'model' : RandomForestClassifier(),
14.        'params' : {
15.            'n_estimators': [130, 140, 145],
16.            'max_depth': [90, 100, 110]
17.        }
18.    },
19.
20.    'xgboost_classifier' : {
21.        'model' : XGBClassifier(verbosity = 0),
22.        'params' : {
23.            'max_depth': [4, 5, 6] ,
24.            'learning_rate': [0.3, 0.4, 0.5],
25.            'subsample': [1, 2],
26.            'gamma' : [5, 6, 7]
27.        }
28.    },
29.
30.    'logistic_regression' : {
31.        'model' : LogisticRegression(solver='liblinear'),
32.        'params' : {
33.            'penalty' : ['l2'],
34.            'C': [0.5, 1.0, 1.5],
35.            'max_iter' : [40, 50, 60]
36.        }
37.    }
38. }

```

## Hyperparameter Tuning

```

1. scores = []
2. best_estimators = {}
3. for algo, mp in model_params.items():
4.     clf = GridSearchCV(mp['model'], mp['params'], cv=5, return_train_score=False)
5.     clf.fit(X_train, y_train)
6.     scores.append({

```

```

7.         'model':algo,
8.         'best_score':clf.best_score_,
9.         'best_params':clf.best_params_
10.    })
11.    best_estimators[algo] = clf.best_estimator_

```

```

1. df = pd.DataFrame(scores, columns=['model', 'best_score', 'best_params'])
2. df

```

	model	best_score	best_params
0	decision_tree	0.956274	{'criterion': 'entropy', 'max_depth': 22, 'min...
1	random_forest	0.974297	{'max_depth': 110, 'n_estimators': 140}
2	xgboost_classifier	0.954612	{'gamma': 5, 'learning_rate': 0.4, 'max_depth'...
3	logistic_regression	0.932720	{'C': 1.5, 'max_iter': 40, 'penalty': 'l2'}

We'll select **XGBClassifier** as a final model.

```

1. clf = best_estimators['xgboost_classifier']
2. accuracy = clf.score(X_test, y_test)
3. prediction = clf.predict(X_test)
4. probs = clf.predict_proba(X_test)[: , 1]
5. roc_auc = roc_auc_score(y_test, probs)
6. print("Accuracy Score :", round(accuracy, 4))
7. print()
8. print("ROU AUC Score :", round(roc_auc, 4))
9. print()
10. print("Classification Report:\n")
11. print(classification_report(y_test, prediction))
12. print()
13. print("Note: As this is an imbalanced dataset accuracy can give us false
    assumptions regarding performance. So, It's better to rely on ROC metricses.")
14. print()
15. fpr, tpr, thresholds = roc_curve(y_test, probs)
16. plt.plot(fpr, tpr, color='darkorange', label='ROC curve (area = %0.2f)' % roc_auc)
17. plt.plot([0, 1], [0, 1], color='navy', linestyle='--')
18. plt.xlim([-0.01, 1.0])
19. plt.ylim([0.0, 1.01])
20. plt.xlabel('False Positive Rate')
21. plt.ylabel('True Positive Rate')
22. plt.title('Receiver operating characteristic')
23. plt.legend(loc="lower right")

```

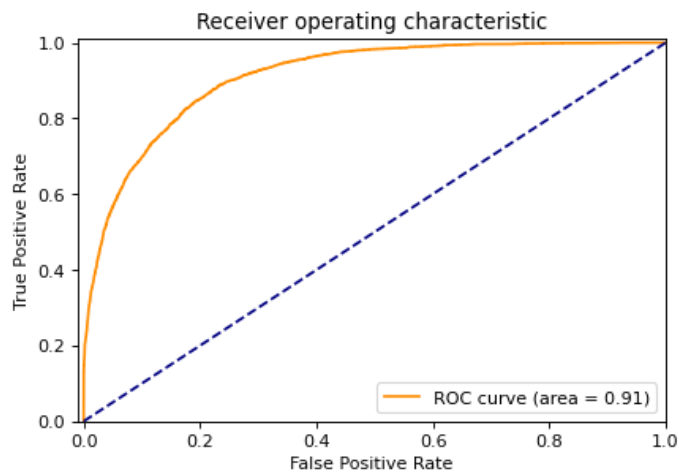
Accuracy Score : 0.8039

ROU AUC Score : 0.9104

Classification Report:

	precision	recall	f1-score	support
0	0.95	0.78	0.86	12435
1	0.55	0.87	0.68	3846
accuracy			0.80	16281
macro avg	0.75	0.83	0.77	16281
weighted avg	0.86	0.80	0.82	16281

Note: As this is an imbalanced dataset accuracy can give us false assumptions regarding performance. So, It's better to rely on ROC metrics.



Now, save this model so that we can use this same model in future.

```
1. pickle.dump(clf, open('final_model', 'wb'))
```