





```
test_data_with_obj = test_data[test_obj_cols]
test_data_with_obj.head()
```

```
workclass  education  marital_status  occupation  relationship  race  sex  native_country  income
0 Private      HS-grad      Married-civ-spouse  Farming-fishing      Husband      White  Male      United-States  <=50K.
1 Local-gov      Assoc-acdm      Married-civ-spouse  Protective-serv      Husband      White  Male      United-States  >50K.
2 Private      Some-college      Married-civ-spouse  Machine-op-inpct      Husband      Black  Male      United-States  >50K.
3 workclass_Not_available      Some-college      Never-married      occupation_Not_available  Own-child      White  Female      United-States  <=50K.
4 Private      9th and Above      Never-married      Other-service      Not-in-family      White  Male      United-States  <=50K.
```

```
In [65]: test_two_classes = ["sex", "income"]
test_multiple_classes = ["workclass", "education", "marital_status", "occupation", "relationship", "race"]

In [66]: from sklearn.preprocessing import LabelEncoder
test_data_two_class = pd.DataFrame()
test_data = LabelEncoder().fit_transform(test_data_with_obj[test_obj_cols],)
test_data = LabelEncoder().fit_transform(test_data_with_obj[test_multiple_classes]),)

In [67]: test_data_multiple_classes = pd.get_dummies(test_data_with_obj[test_multiple_classes])
test_data_multiple_classes.head()
```

```
workclass_Federal-gov  workclass_Local-gov  workclass_Never-worked  workclass_Private  workclass_Self-emp-inc  workclass_Self-emp-not-inc  workclass_State-gov  workclass_Without-pay  workclass_Not_avail
0 0 0 0 1 0 0 0 0 0
1 0 1 0 0 0 0 0 0 0
2 0 0 0 0 1 0 0 0 0
3 0 0 0 0 0 0 0 0 1
4 0 0 0 1 0 0 0 0 0
5 rows x 9 columns
```

```
In [68]: test_list = [test_data_two_class, test_data_multiple_classes]
test_final_obj_data = pd.concat(test_list, axis = 1)
test_final_obj_data.head()
```

```
sex  income  workclass_Federal-gov  workclass_Local-gov  workclass_Never-worked  workclass_Private  workclass_Self-emp-inc  workclass_Self-emp-not-inc  workclass_State-gov  workclass_Without-pay  workclass_Not_avail
0 1 0 0 0 0 0 0 0 0 0 0
1 0 1 0 0 0 0 0 0 0 0 0
2 0 1 1 0 0 0 0 0 0 0 0
3 0 0 0 0 0 0 0 0 0 0 0
4 1 0 0 0 0 0 1 0 0 0 0
5 rows x 11 columns
```

```
In [69]: test_data_with_num_hours_per_week = test_data_hours_per_week[100]
test_data_with_num_age = test_data_with_num_age[100]
test_data_with_num_education_num = test_data_with_num_education_num[100]
```

C:\ProgramData\Anaconda3\lib\site-packages\pandas\core\generic.py:5096: SettingWithCopyWarning:
A value is being stored to test\_data\_with\_num by a slice from a DataFrame.
Try using .loc[row\_indexer,col\_indexer] = value instead

See the caveats in the documentation: http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing-view-versus-copy
self[name] = value

```
In [70]: from sklearn.preprocessing import MinMaxScaler
sc = MinMaxScaler()
test_list_scaler_values = ["nlwgt", "capital_gain", "capital_loss"]
test_data_array = sc.fit_transform(test_data_with_num[test_list_scaler_values])
test_data_array[:5]
```

```
array([[0.05167688, 0.      , 0.      ],
       [0.21901083, 0.      , 0.      ],
       [0.0094186 , 0.07686975, 0.      ],
       [0.06094151, 0.      , 0.      ],
       [0.12539773, 0.      , 0.      ]])
```

```
In [71]: test_column_values = ["nlwgt", "capital_gain", "capital_loss"]
test_df = pd.DataFrame(data = test_data_array,
                        columns = test_column_values)
test_df.tail()
```

```
nlwgt  capital_gain  capital_loss
16275  0.136723  0.000000  0.0
16276  0.208484  0.000000  0.0
16277  0.244762  0.000000  0.0
16278  0.047866  0.054551  0.0
16279  0.114195  0.000000  0.0
```

```
In [72]: test_list_num = [test_data_with_num[["sex", "hours_per_week", "education_num"]], test_df]
test_list_num_data = pd.concat(test_list_num, axis = 1)
```

```
In [73]: test_final_df_list = [test_final_num_data, test_final_obj_data]
test_final_data = pd.concat(test_final_df_list, axis = 1)
test_final_data.head()
```

```
age  hours_per_week  education_num  nlwgt  capital_gain  capital_loss  sex  income  workclass_Federal-gov  workclass_Local-gov  ...  n
0  0.38  0.5  0.09  0.051677  0.000000  0.0  1  0  0  0  0  ...  0
1  0.28  0.4  0.12  0.219011  0.000000  0.0  1  1  0  1  1  ...  0
2  0.44  0.4  0.10  0.009418  0.076869  0.0  1  1  0  0  0  ...  0
3  0.18  0.3  0.10  0.060942  0.000000  0.0  0  0  0  0  0  ...  0
4  0.34  0.3  0.06  0.125398  0.000000  0.0  1  0  0  0  0  ...  1
5 rows x 14 columns
```

```
In [74]: test_X = test_final_data.drop("income", axis = 1)
test_Y = test_final_data.income
```

```
In [75]: from sklearn.decomposition import PCA
pca = PCA(n_components= 5)
test_transformed_data = pca.fit_transform(test_X)
test_transformed_data
```

```
array([[ -0.99487154,  0.75211029, -0.11469905],
       [-1.04730551, -0.53059653,  0.17502619],
       [-0.79461052, -0.05445943, -0.36029345],
       ...,
       [-1.01022029, -0.33803172, -0.04831659],
       [ 0.48221337, -0.02828259,  0.02320211],
       [-1.14022289, -0.75833189,  0.23749085]])
```

## Using Training Model For Testing

```
In [77]: from sklearn.tree import DecisionTreeClassifier
dec = DecisionTreeClassifier(criterion='entropy', max_depth=10)
```

```
DecisionTreeClassifier(class_weight=None, criterion='entropy', max_depth=10,
max_features=None, max_leaf_nodes=None,
min_impurity_decrease=0.0, min_impurity_split=None,
min_samples_leaf=1, min_samples_split=2,
min_weight_fraction_leaf=0.0, presort=False,
random_state=None, splitter='best')
```

```
In [80]: dec.fit(transformed_data, Y)
```

```
DecisionTreeClassifier(class_weight=None, criterion='entropy', max_depth=10,
max_features=None, max_leaf_nodes=None,
min_impurity_decrease=0.0, min_impurity_split=None,
min_samples_leaf=1, min_samples_split=2,
min_weight_fraction_leaf=0.0, presort=False,
random_state=None, splitter='best')
```

```
In [81]: Y_pred = dec.predict(test_transformed_data)
```

```
In [82]: (test_Y, Y_pred)
```

```
(16280, 16280)
```

```
In [83]: MSE = np.square(np.subtract(test_Y,Y_pred)).mean()
MSE
```

```
0.2516584760584767
```

```
In [84]: list_pred = []
for i in Y_pred:
    if i == 1:
        list_pred.append("<=50")
    else:
        list_pred.append(">50")
```

```
In [85]: list_actual = []
for i in test_Y:
    if i == 1:
        list_actual.append("<=50")
    else:
        list_actual.append(">50")
```

```
In [86]: final_pred = pd.DataFrame({
    "Y_actual": list_actual,
    "Y_pred": list_pred
})
```

```
In [87]: final_pred.head()
```

```
Y_actual  Y_pred
0 <=50      <=50
1 >50       >50
2 >50      <=50
3 <=50      <=50
4 <=50      <=50
```

```
In [ ]:
```