

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 from sklearn.datasets.base import Bunch
5 from sklearn.preprocessing import LabelEncoder
6 from sklearn.base import BaseEstimator, TransformerMixin
```

```
↳ /usr/local/lib/python3.6/dist-packages/statsmodels/tools/_testing.py:19: FutureWarning: pandas.util.testing is deprecated
    import pandas.util.testing as tm
/usr/local/lib/python3.6/dist-packages/sklearn/utils/deprecation.py:144: FutureWarning: The sklearn.datasets.base module is deprecated
    warnings.warn(message, FutureWarning)
```

▼ Read in the data and explore it

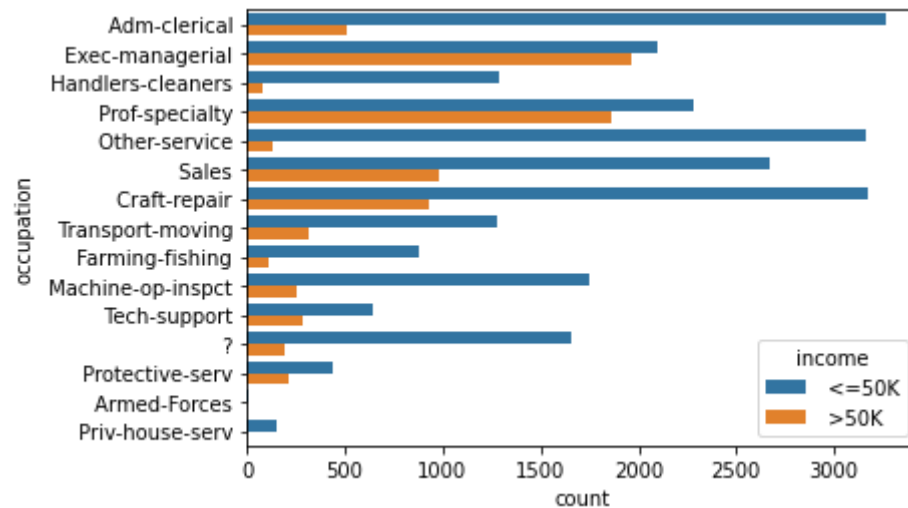
```
1
2 on-num', 'marital-status', 'occupation', 'relationship', 'race', 'sex', 'capital-gain', 'capital-loss', 'hours-per-week', 'native-country']
```

```
1 #reading in the data
2 adult_data_train=pd.read_csv('adult.data',names=col_names)
3 adult_data_predict= pd.read_csv('adult.test', names=col_names, skiprows=1)
```

```
1 sns.countplot(y='occupation', hue='income', data=adult_data_train,)
```

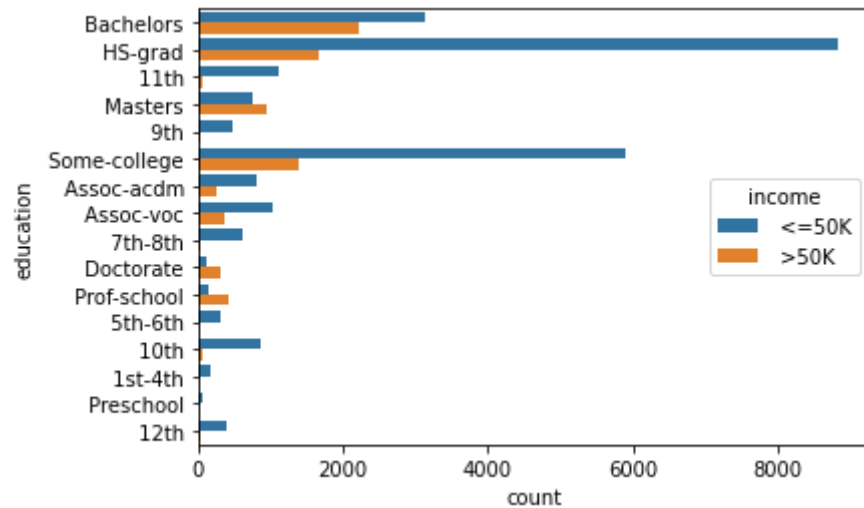
```
↳
```

<matplotlib.axes._subplots.AxesSubplot at 0x7f2a2c813b00>



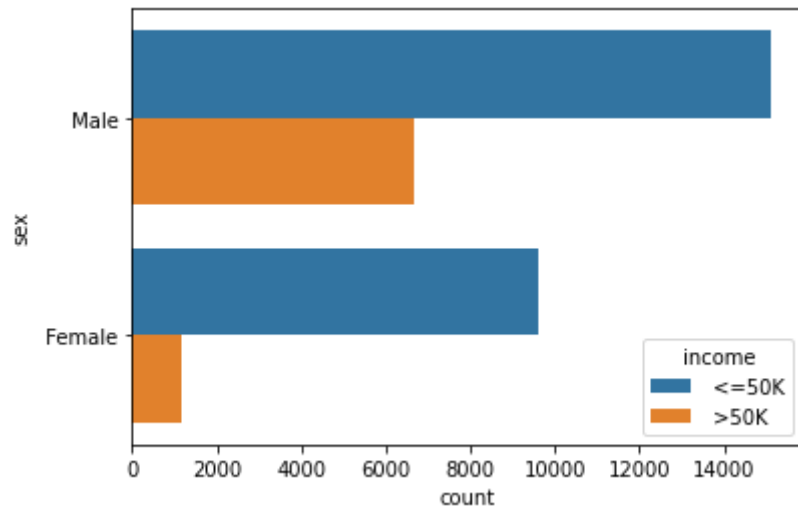
```
1 sns.countplot(y='education', hue='income', data=adult_data_train,)
```

↳ <matplotlib.axes._subplots.AxesSubplot at 0x7f2a2c6c6e80>



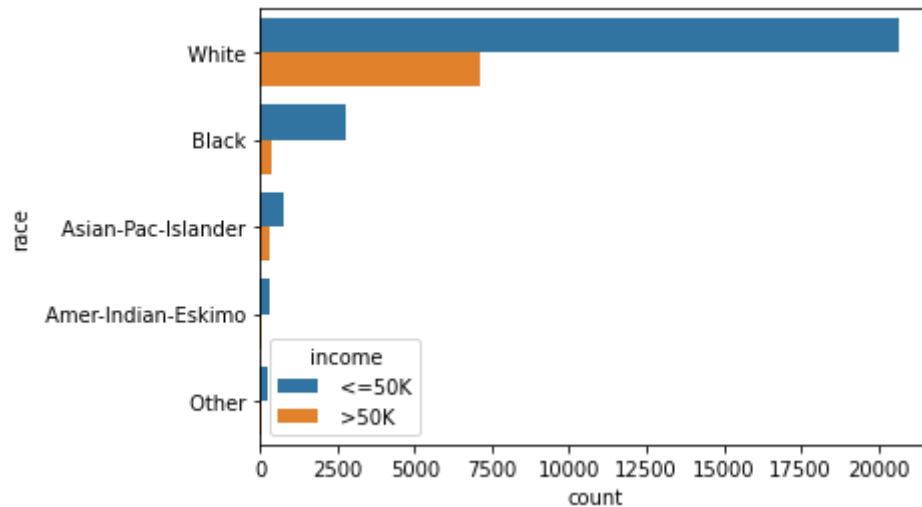
```
1 sns.countplot(y='sex', hue='income', data=adult_data_train,)
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f2a2c1ce048>
```



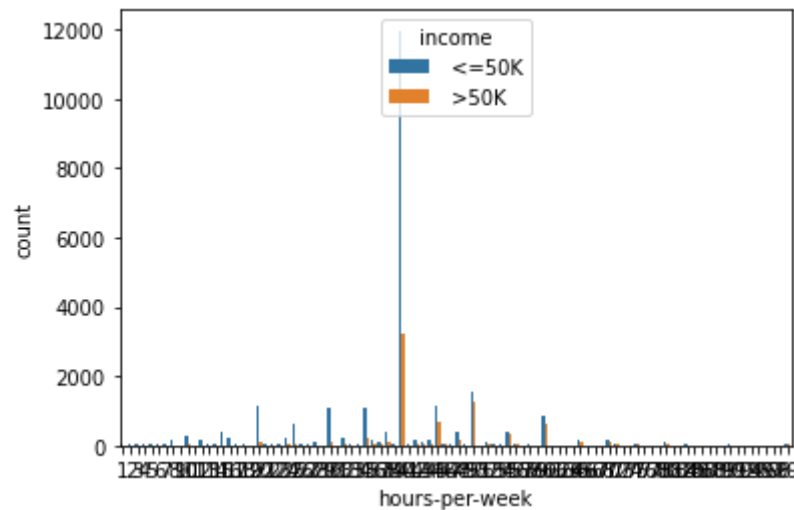
```
1 sns.countplot(y='race', hue='income', data=adult_data_train,)
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f2a2c138a90>
```



```
1 sns.countplot(x='hours-per-week', hue='income', data=adult_data_train,)
```

↳ <matplotlib.axes._subplots.AxesSubplot at 0x7f2a2c0796d8>



▼ Preprocessing the data

```
1 #1. Replace the Nan/Missing data values with mode values
2 attr_list=list(col_names)
3 attr_list.remove('income')
4
5 for col in attr_list:
6     mode=adult_data_train[col].mode()
7     adult_data_train[col].fillna(mode)
8
9 for col in attr_list:
10     mode=adult_data_predict[col].mode()
11     adult_data_predict[col].fillna(mode)
```

```
1 #2. For encoding the categorical data
2 class EncodeCategorical(BaseEstimator, TransformerMixin):
3     """
4     Encodes a specified list of columns or all columns if None
```

```

4     Encodes a specified list of columns or all columns if None.
5     """
6     def __init__(self, columns=None):
7         self.columns = columns
8         self.encoders = None
9     def fit(self, data, target=None):
10        """
11        Expects a data frame with named columns to encode.
12        """
13        # Encode all columns if columns is None
14        if self.columns is None:
15            self.columns = data.columns
16        # Fit a label encoder for each column in the data frame
17        self.encoders = {
18            column: LabelEncoder().fit(data[column])
19            for column in self.columns
20        }
21        return self
22    def transform(self, data):
23        """
24        Uses the encoders to transform a data frame.
25        """
26        output = data.copy()
27        for column, encoder in self.encoders.items():
28            output[column] = encoder.transform(data[column])
29        return output
30
31 encoder = EncodeCategorical(col_names)
32 data = encoder.fit_transform(adult_data_train)
33 data_predict=encoder.fit_transform(adult_data_predict)

```

```

1 #3. Split into the train and test sets
2 from sklearn.model_selection import train_test_split
3 X_train,X_test,y_train,y_test=train_test_split(data,adult_data_train['income'],test_size=0.2)

```

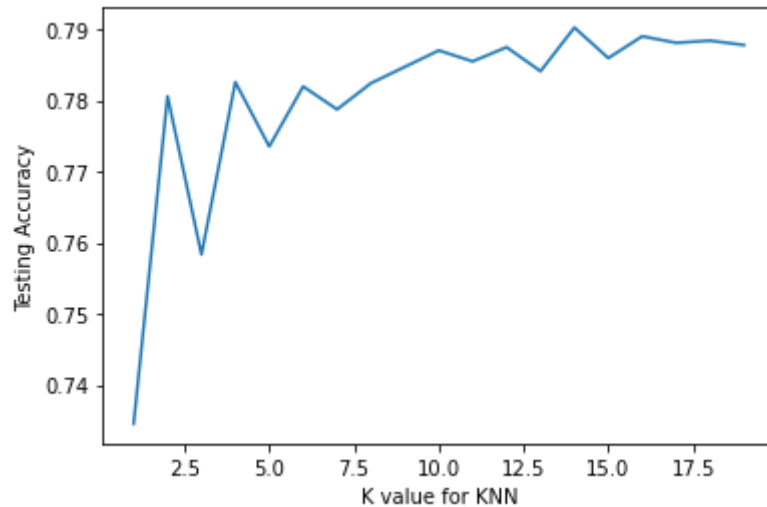
The KNN Model

```

1 #1.Find the optimum value of k
2 from sklearn.neighbors import KNeighborsClassifier
3 from sklearn import metrics
4
5 k_range=range(1,20)
6 scores={}
7 scores_list=[]
8
9 for k in k_range:
10     knn=KNeighborsClassifier(n_neighbors=k)
11     knn.fit(X_train,y_train)
12     y_pred=knn.predict(X_test)
13     sc=metrics.accuracy_score(y_test,y_pred)
14     scores[k]=sc
15     scores_list.append(sc)
16
17 plt.plot(k_range,scores_list)
18 plt.xlabel('K value for KNN')
19 plt.ylabel('Testing Accuracy')

```

☞ Text(0, 0.5, 'Testing Accuracy')



```

1 #Train the knn model with optimum value of k

```

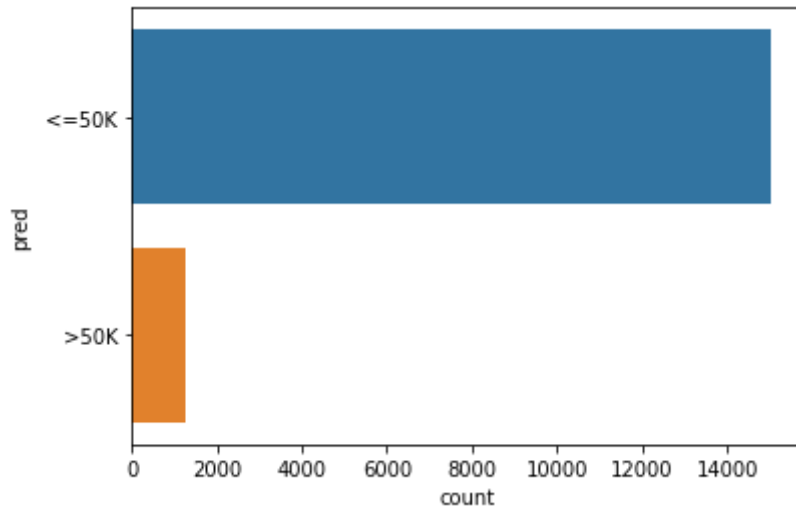
```
2 k_optimum=17
3 knn=KNeighborsClassifier(n_neighbors=k_optimum)
4 knn.fit(X_train,y_train)
5 knn.score(X_test, y_test)
```

0.7881160755412252

```
1 #predict the value for the test sets
2 prediction=knn.predict(data_predict)
```

```
1 pred=pd.DataFrame(prediction,columns=['pred'])
2 sns.countplot(y='pred',data=pred,)
```

<matplotlib.axes._subplots.AxesSubplot at 0x7f2a29551e10>



1

