In [5]:
```python
import numpy as np
import pandas as pd
%matplotlib inline
import matplotlib.pyplot as plt
import seaborn as sns
```
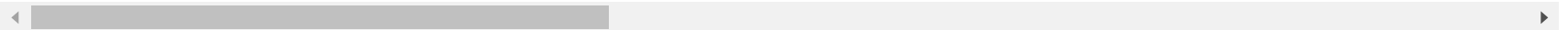
In [20]:
```python
dataset = pd.read_csv("HR_Employee_Attrition_IBM.csv")
dataset.head()
```

Out[20]:

| | Age | Attrition | BusinessTravel | DailyRate | Department | DistanceFromHome | Education | EducationField | EmployeeCount | EmployeeN |
|---|-----|-----------|----------------|-----------|------------|------------------|-----------|----------------|---------------|-----------|
| 0 | 41 | Yes | Travel_Rarely | 1102 | Sales | 1 | 2 | Life Sciences | 1 | |
| 1 | 49 | No | Travel_Frequently | 279 | Research & Development | 8 | 1 | Life Sciences | 1 | |
| 2 | 37 | Yes | Travel_Rarely | 1373 | Research & Development | 2 | 2 | Other | 1 | |
| 3 | 33 | No | Travel_Frequently | 1392 | Research & Development | 3 | 4 | Life Sciences | 1 | |
| 4 | 27 | No | Travel_Rarely | 591 | Research & Development | 2 | 1 | Medical | 1 | |

5 rows × 35 columns

In [9]: ```
dataset.info()
```

```
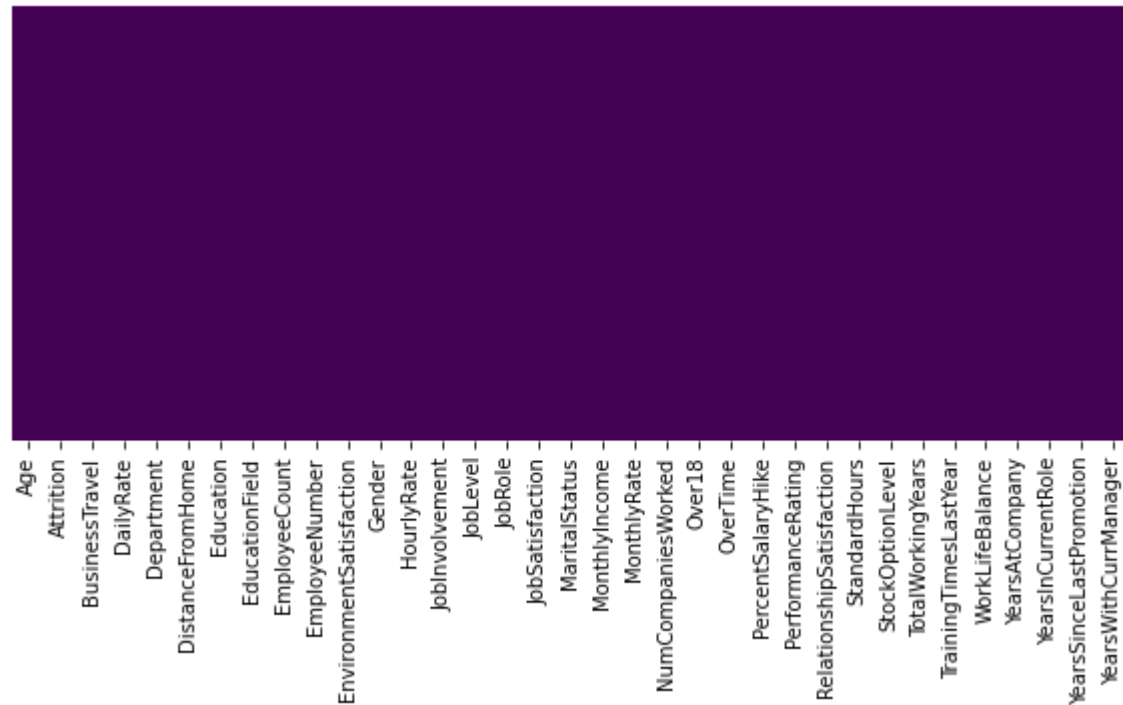<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1470 entries, 0 to 1469
Data columns (total 35 columns):
 #   Column                    Non-Null Count  Dtype
---  ------                    --------------  -----
 0   Age                       1470 non-null   int64
 1   Attrition                 1470 non-null   object
 2   BusinessTravel            1470 non-null   object
 3   DailyRate                 1470 non-null   int64
 4   Department                1470 non-null   object
 5   DistanceFromHome          1470 non-null   int64
 6   Education                 1470 non-null   int64
 7   EducationField            1470 non-null   object
 8   EmployeeCount             1470 non-null   int64
 9   EmployeeNumber            1470 non-null   int64
 10  EnvironmentSatisfaction   1470 non-null   int64
 11  Gender                    1470 non-null   object
 12  HourlyRate                1470 non-null   int64
 13  JobInvolvement            1470 non-null   int64
 14  JobLevel                  1470 non-null   int64
 15  JobRole                   1470 non-null   object
 16  JobSatisfaction           1470 non-null   int64
 17  MaritalStatus             1470 non-null   object
 18  MonthlyIncome             1470 non-null   int64
 19  MonthlyRate               1470 non-null   int64
 20  NumCompaniesWorked        1470 non-null   int64
 21  Over18                    1470 non-null   object
 22  OverTime                  1470 non-null   object
 23  PercentSalaryHike         1470 non-null   int64
 24  PerformanceRating         1470 non-null   int64
 25  RelationshipSatisfaction  1470 non-null   int64
 26  StandardHours             1470 non-null   int64
 27  StockOptionLevel          1470 non-null   int64
 28  TotalWorkingYears         1470 non-null   int64
 29  TrainingTimesLastYear     1470 non-null   int64
 30  WorkLifeBalance           1470 non-null   int64
 31  YearsAtCompany            1470 non-null   int64
 32  YearsInCurrentRole        1470 non-null   int64
 33  YearsSinceLastPromotion   1470 non-null   int64
 34  YearsWithCurrManager      1470 non-null   int64
dtypes: int64(26), object(9)
memory usage: 402.1+ KB
```

heatmap to check the missing value

In [10]: 
```
plt.figure(figsize =(10, 4))
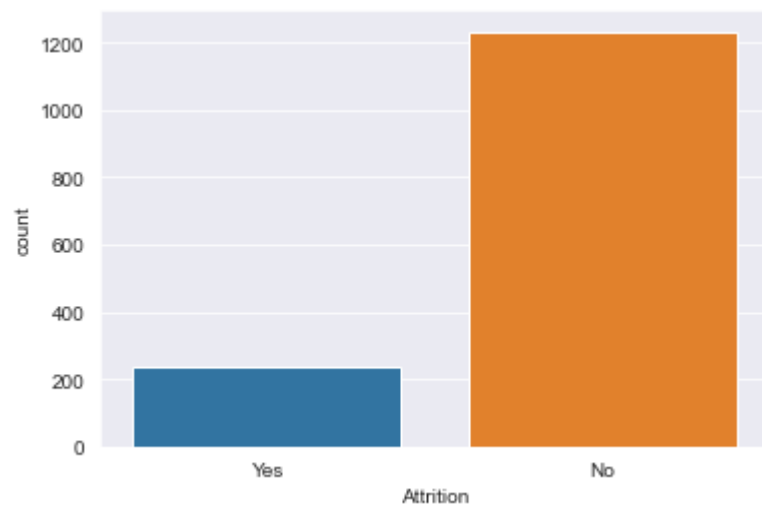sns.heatmap(dataset.isnull(), yticklabels = False, cbar = False, cmap ='viridis')
```

Out[10]: <matplotlib.axes._subplots.AxesSubplot at 0x27adf548e80>



So, we can see that there are no missing values in the dataset.

In [11]:
```python
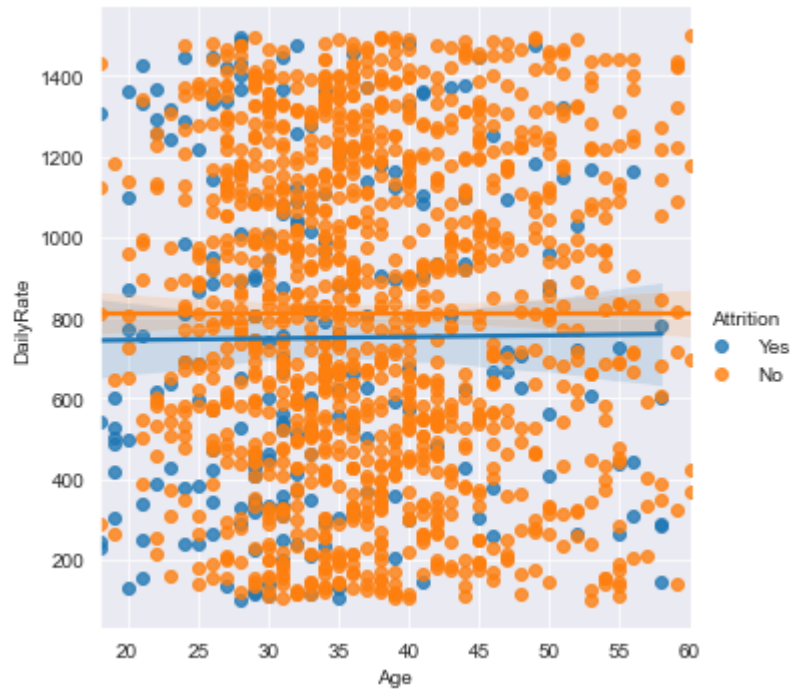sns.set_style('darkgrid')
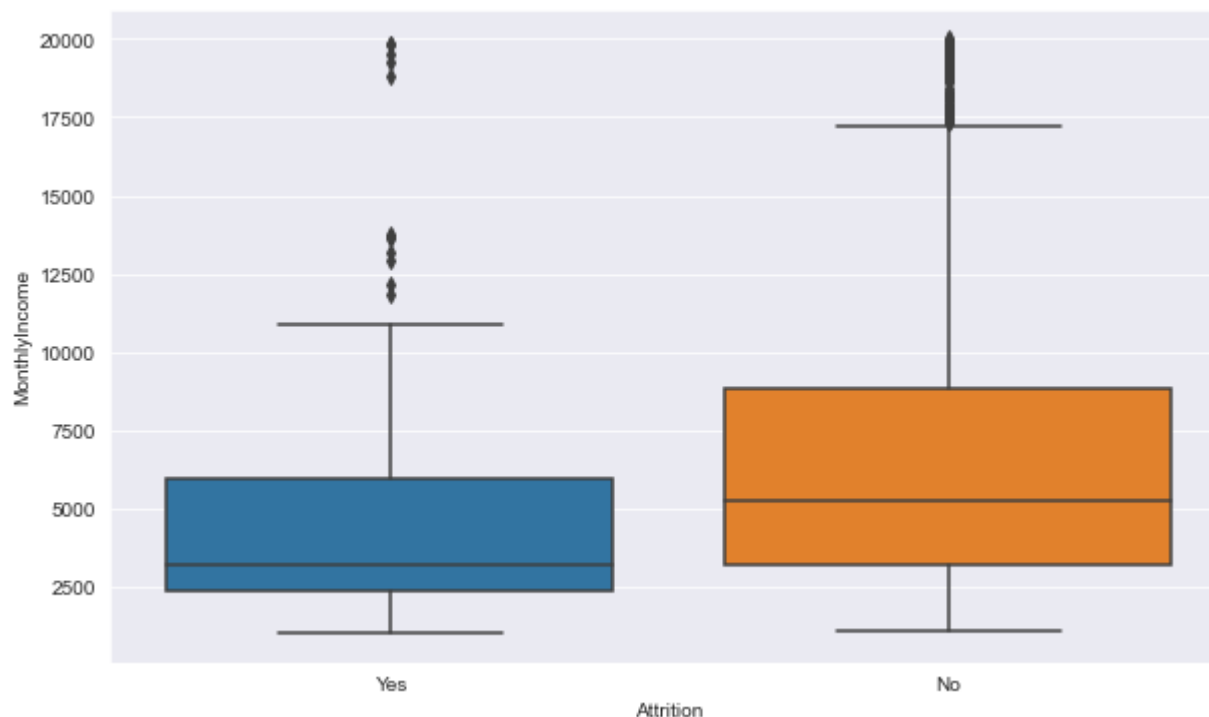sns.countplot(x ='Attrition', data = dataset)
```

Out[11]: <matplotlib.axes._subplots.AxesSubplot at 0x27adf5f98b0>

In [12]: `sns.lmplot(x = 'Age', y = 'DailyRate', hue = 'Attrition', data = dataset)`

Out[12]: `<seaborn.axisgrid.FacetGrid at 0x27adf5d68b0>`

In [13]:
```python
plt.figure(figsize =(10, 6))
sns.boxplot(y ='MonthlyIncome', x ='Attrition', data = dataset)
```

Out[13]: `<matplotlib.axes._subplots.AxesSubplot at 0x27adff9c640>`



In the dataset there are 4 irrelevant columns, i.e:EmployeeCount, EmployeeNumber, Over18 and StandardHour. So, we have to remove these for more accuracy.

In [21]:
```python
dataset.drop('EmployeeCount', axis = 1, inplace = True)
dataset.drop('StandardHours', axis = 1, inplace = True)
dataset.drop('EmployeeNumber', axis = 1, inplace = True)
dataset.drop('Over18', axis = 1, inplace = True)
print(dataset.shape)
```

(1470, 31)

In [22]:
```python
y = dataset.iloc[:, 1]
X = dataset
X.drop('Attrition', axis = 1, inplace = True)
```

```
In [23]: from sklearn.preprocessing import LabelEncoder
         lb = LabelEncoder()
         y = lb.fit_transform(y)
```

In the dataset there are 7 categorical data, so we have to change them to int data, i.e we hava to create 7 dummy variable for better accuracy.

```
In [25]: dum_BusinessTravel = pd.get_dummies(dataset['BusinessTravel'],prefix ='BusinessTravel')
         dum_Department = pd.get_dummies(dataset['Department'],prefix ='Department')
         dum_EducationField = pd.get_dummies(dataset['EducationField'],prefix ='EducationField')
         dum_Gender = pd.get_dummies(dataset['Gender'], prefix ='Gender', drop_first = True)
         dum_JobRole = pd.get_dummies(dataset['JobRole'],prefix ='JobRole')
         dum_MaritalStatus = pd.get_dummies(dataset['MaritalStatus'],prefix ='MaritalStatus')
         dum_OverTime = pd.get_dummies(dataset['OverTime'], prefix ='OverTime', drop_first = True)
         # Adding these dummy variable to input X
         X = pd.concat([X, dum_BusinessTravel, dum_Department,dum_EducationField, dum_Gender, dum_JobRole,dum_MaritalS
         tatus, dum_OverTime], axis = 1)
         # Removing the categorical data
         X.drop(['BusinessTravel', 'Department', 'EducationField', 'Gender', 'JobRole', 'MaritalStatus', 'OverTime'],
         axis = 1, inplace = True)

         print(X.shape)
         print(y.shape)
```

```
(1470, 49)
(1470,)
```

```
In [26]: from sklearn.model_selection import train_test_split
         X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25, random_state = 40)
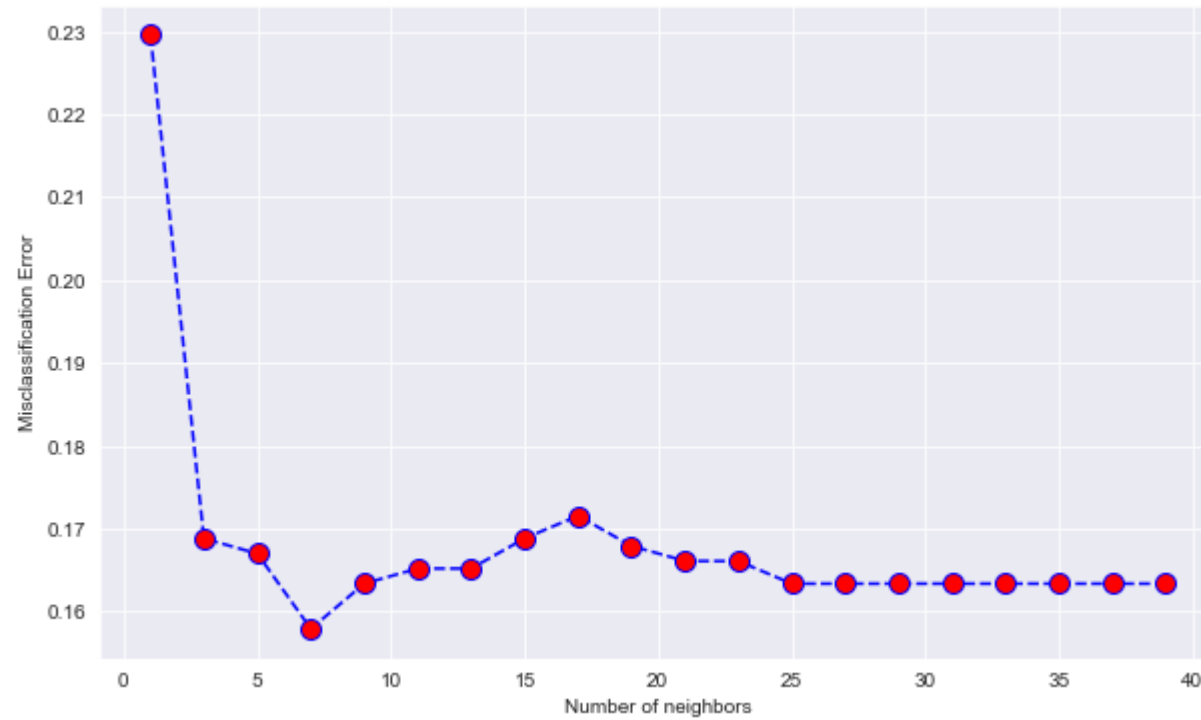```

Apply KNN Algo

In [27]:
```python
from sklearn.neighbors import KNeighborsClassifier
neighbors = []
cv_scores = []

from sklearn.model_selection import cross_val_score
# perform 10 fold cross validation
for k in range(1, 40, 2):
    neighbors.append(k)
    knn = KNeighborsClassifier(n_neighbors = k)
    scores = cross_val_score(knn, X_train, y_train, cv = 10, scoring = 'accuracy')
    cv_scores.append(scores.mean())
error_rate = [1-x for x in cv_scores]

# determining the best k
optimal_k = neighbors[error_rate.index(min(error_rate))]
print('The optimal number of neighbors is % d ' % optimal_k)

# plot misclassification error versus k
plt.figure(figsize = (10, 6))
plt.plot(range(1, 40, 2), error_rate, color ='blue', linestyle ='dashed', marker ='o',markerfacecolor ='red',
markersize = 10)
plt.xlabel('Number of neighbors')
plt.ylabel('Misclassification Error')
plt.show()
```

The optimal number of neighbors is  7



The optimal number of neighbors is 7

In [28]:
```python
from sklearn.model_selection import cross_val_predict, cross_val_score
from sklearn.metrics import accuracy_score, classification_report
from sklearn.metrics import confusion_matrix

def print_score(clf, X_train, y_train, X_test, y_test, train = True):
    if train:
        print("Train Result:")
        print("------------")
        print("Classification Report: \n {}\n".format(classification_report(
                y_train, clf.predict(X_train))))
        print("Confusion Matrix: \n {}\n".format(confusion_matrix(
                y_train, clf.predict(X_train))))

        res = cross_val_score(clf, X_train, y_train,cv = 10, scoring ='accuracy')
        print("Average Accuracy: \t {0:.4f}".format(np.mean(res)))
        print("Accuracy SD: \t\t {0:.4f}".format(np.std(res)))
        print("accuracy score: {0:.4f}\n".format(accuracy_score(y_train, clf.predict(X_train))))
        print("----------------------------------------------------------")


    elif train == False:
        print("Test Result:")
        print("-----------")
        print("Classification Report: \n {}\n".format(classification_report(y_test, clf.predict(X_test))))
        print("Confusion Matrix: \n {}\n".format(confusion_matrix(y_test, clf.predict(X_test))))
        print("accuracy score: {0:.4f}\n".format(accuracy_score(y_test, clf.predict(X_test))))
        print("----------------------------------------------------------")

knn = KNeighborsClassifier(n_neighbors = 7)
knn.fit(X_train, y_train)
print_score(knn, X_train, y_train, X_test, y_test, train = True)
print_score(knn, X_train, y_train, X_test, y_test, train = False)
```

```
Train Result:
------------
Classification Report:
              precision    recall  f1-score   support

           0       0.86      0.99      0.92       922
           1       0.83      0.19      0.32       180

    accuracy                           0.86      1102
   macro avg       0.85      0.59      0.62      1102
weighted avg       0.86      0.86      0.82      1102


Confusion Matrix:
 [[915    7]
 [145   35]]

Average Accuracy:       0.8421
Accuracy SD:            0.0148
accuracy score: 0.8621


-------------------------------------------------------------
Test Result:
-----------
Classification Report:
              precision    recall  f1-score   support

           0       0.84      0.96      0.90       311
           1       0.14      0.04      0.06        57

    accuracy                           0.82       368
   macro avg       0.49      0.50      0.48       368
weighted avg       0.74      0.82      0.77       368


Confusion Matrix:
 [[299  12]
 [ 55   2]]

accuracy score: 0.8179

-------------------------------------------------------------
```

In [ ]: