

Exercise : descriptive summary statistics

Dataset : EmployeeAttrition.csv

```
In [1]: import pandas
import matplotlib.pyplot as plot
```

creating a dataframe from the csv file

```
In [2]: df = pandas.read_csv("Y:\DA LAB\EmployeeAttrition.csv")
```

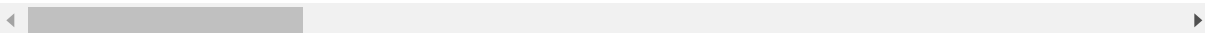
creating a summary statistics using describe()

```
In [3]: df.describe(include='all')
```

Out[3]:

	Age	Attrition	BusinessTravel	DailyRate	Department	DistanceFrom
<b>count</b>	1470.000000	1470	1470	1470.000000	1470	1470.000000
<b>unique</b>	NaN	2	3	NaN	3	NaN
<b>top</b>	NaN	No	Travel_Rarely	NaN	Research & Development	NaN
<b>freq</b>	NaN	1233	1043	NaN	961	NaN
<b>mean</b>	36.923810	NaN	NaN	802.485714	NaN	9.192517
<b>std</b>	9.135373	NaN	NaN	403.509100	NaN	8.106864
<b>min</b>	18.000000	NaN	NaN	102.000000	NaN	1.000000
<b>25%</b>	30.000000	NaN	NaN	465.000000	NaN	2.000000
<b>50%</b>	36.000000	NaN	NaN	802.000000	NaN	7.000000
<b>75%</b>	43.000000	NaN	NaN	1157.000000	NaN	14.000000
<b>max</b>	60.000000	NaN	NaN	1499.000000	NaN	29.000000

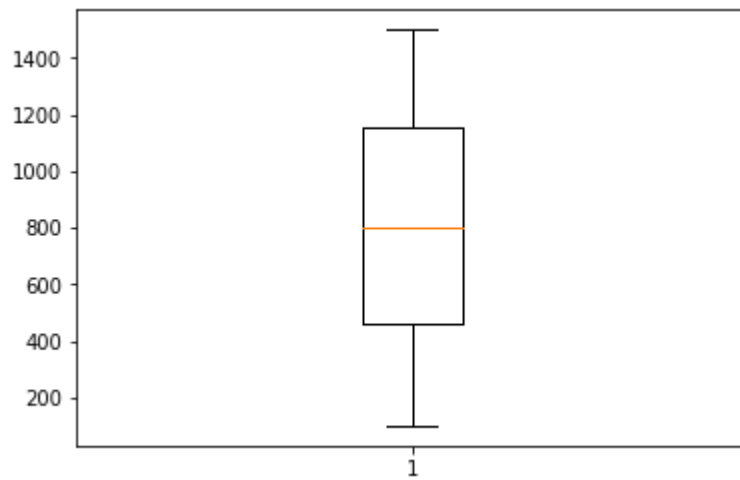
11 rows × 7 columns



creating a boxplot of all numeric attributes using pandas boxplot()



```
In [6]: plot.boxplot(df['DailyRate'])  
plot.show()
```



Checking if there are any missing values

```
In [7]: df.isnull().any()
```

```
Out[7]: Age                False
Attrition                 False
BusinessTravel            False
DailyRate                 False
Department                False
DistanceFromHome          False
Education                 False
EducationField             False
EmployeeCount             False
EmployeeNumber            False
EnvironmentSatisfaction    False
Gender                    False
HourlyRate                 False
JobInvolvement            False
JobLevel                  False
JobRole                   False
JobSatisfaction            False
MaritalStatus             False
MonthlyIncome             False
MonthlyRate               False
NumCompaniesWorked        False
Over18                    False
OverTime                  False
PercentSalaryHike         False
PerformanceRating         False
RelationshipSatisfaction   False
StandardHours             False
StockOptionLevel          False
TotalWorkingYears         False
TrainingTimesLastYear     False
WorkLifeBalance           False
YearsAtCompany            False
YearsInCurrentRole        False
YearsSinceLastPromotion    False
YearsWithCurrManager      False
dtype: bool
```

Creating a scatter plot of Age vs Attrition

```
In [9]: plot.scatter(df['Age'],df['Attrition'])  
plot.show()
```

