

```
In [4]: import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
filelocation="C:\Users\A674682\Desktop\workknot\ipynbfiles\haberman.csv"
filelocation1="C:\Users\akars_w3k1ua\Downloads\haberman.csv"

In [5]: #data=pd.read_csv(filelocation1,names=["Age","Year","axilnodes","survival"])
#data[:9]

Out[5]:
```

	Age	Year	axilnodes	survival
0	30	64	1	1
1	30	62	3	1
2	30	65	0	1
3	31	59	2	1
4	31	65	4	1
5	33	58	10	1
6	33	60	0	1
7	34	59	0	2
8	34	66	9	2

```
In [40]: print(data.columns)
print(data.shape)
print(data['survival'].value_counts())

Index(['Age', 'Year', 'axilnodes', 'survival'], dtype='object')
(306, 4)
1    225
2     81
Name: survival, dtype: int64
```

## HIGH LEVEL STATISTICS

- 1.The Data consists of 305 rows or datapoints and 4 columns or features
- 2 There are 4 features/variables out of which 3 are independent variables and one named survival which is the last column is dependent variable or class label
- 3.The columns being Age i.e Age of operation, Year -Year of operation, axilnodes- no of +ve auxiliary nodes and survival - sttus of how long the patient survived
- 4.The data points of survival does not seem to be balanced

## OBJECTIVE

**Classifing wether the patient survives 5 years or more based on age,year,auxillary nodes**

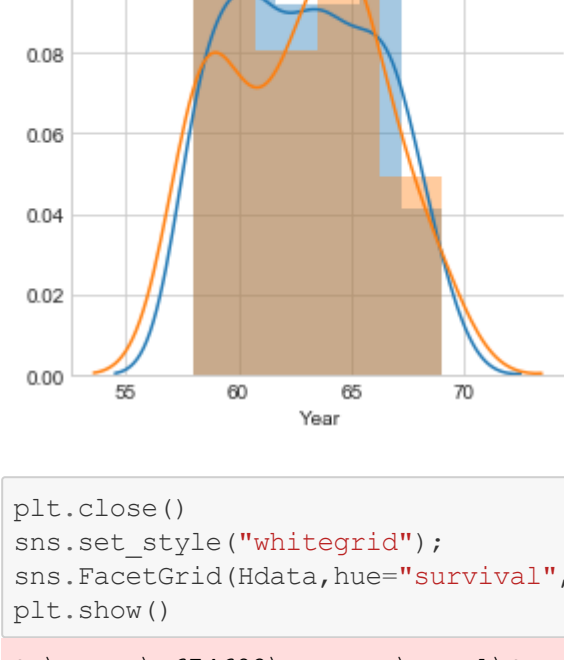
### Univariate Analysis:-

**Histogram of the Independent features :-**

```
In [6]: #data_1=#data[data["survival"]==1]
#data_2=#data[data["survival"]==2]
plt.close()
sns.set_style("whitegrid");
sns.FacetGrid(data,hue="survival",size=4).map(sns.distplot,'Age').add_legend()
plt.show()
```

C:\Users\akars\_w3k1ua\Anaconda3\lib\site-packages\matplotlib\axes\\_axes.py:6462: UserWarning: The 'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.  
warnings.warn("The 'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.", DeprecationWarning)

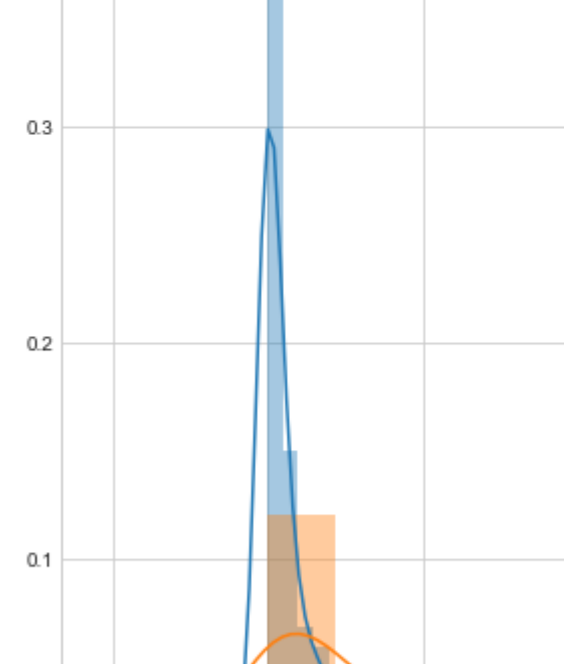
C:\Users\akars\_w3k1ua\Anaconda3\lib\site-packages\matplotlib\axes\\_axes.py:6462: UserWarning: The 'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.  
warnings.warn("The 'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.", DeprecationWarning)



```
In [50]: plt.close()
sns.set_style("whitegrid");
sns.FacetGrid(data,hue="survival",size=9).map(sns.distplot,'Year').add_legend()
plt.show()
```

C:\Users\A674682\AppData\Local\Continuum\anaconda3\lib\site-packages\matplotlib\axes\\_axes.py:6462: UserWarning: The 'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.  
warnings.warn("The 'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.", DeprecationWarning)

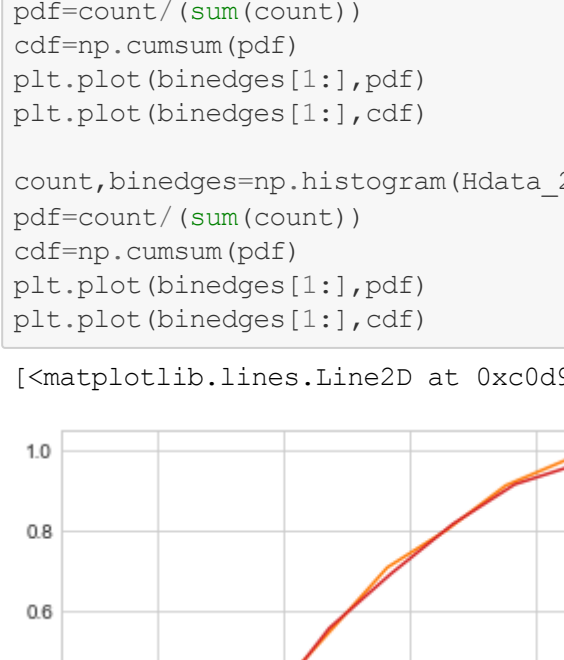
C:\Users\A674682\AppData\Local\Continuum\anaconda3\lib\site-packages\matplotlib\axes\\_axes.py:6462: UserWarning: The 'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.  
warnings.warn("The 'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.", DeprecationWarning)



```
In [56]: plt.close()
sns.set_style("whitegrid");
sns.FacetGrid(data,hue="survival",size=9).map(sns.distplot,'axilnodes').add_legend()
plt.show()
```

C:\Users\A674682\AppData\Local\Continuum\anaconda3\lib\site-packages\matplotlib\axes\\_axes.py:6462: UserWarning: The 'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.  
warnings.warn("The 'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.", DeprecationWarning)

C:\Users\A674682\AppData\Local\Continuum\anaconda3\lib\site-packages\matplotlib\axes\\_axes.py:6462: UserWarning: The 'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.  
warnings.warn("The 'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.", DeprecationWarning)



### Observation :-

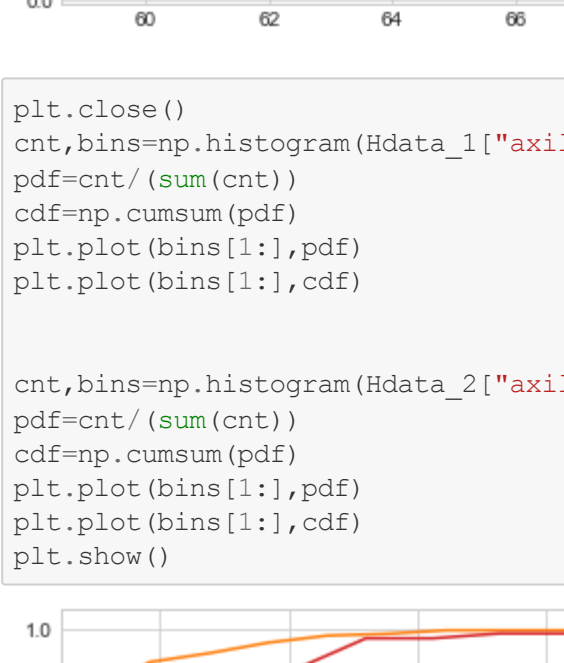
- 1) We can see that it is not possible to classify wether the patient survives more than 5 years or less than five yers by looking into each independent variables as Histogram of each feature there is an overlap of the survival status 1 & 2 and hence Single feature alone does not seem useful for classification

### PDF and CDF

```
In [65]: count,binedges=np.histogram(data_1["Age"],bins=10,density=True)
pdf=count/(sum(count))
cdf=np.cumsum(pdf)
plt.plot(binedges[1:],pdf)
plt.plot(binedges[1:],cdf)

count,binedges=np.histogram(data_2["Age"],bins=10,density=True)
pdf=count/(sum(count))
cdf=np.cumsum(pdf)
plt.plot(binedges[1:],pdf)
plt.plot(binedges[1:],cdf)

Out[65]: <matplotlib.lines.Line2D at 0xc0d9a90>
```



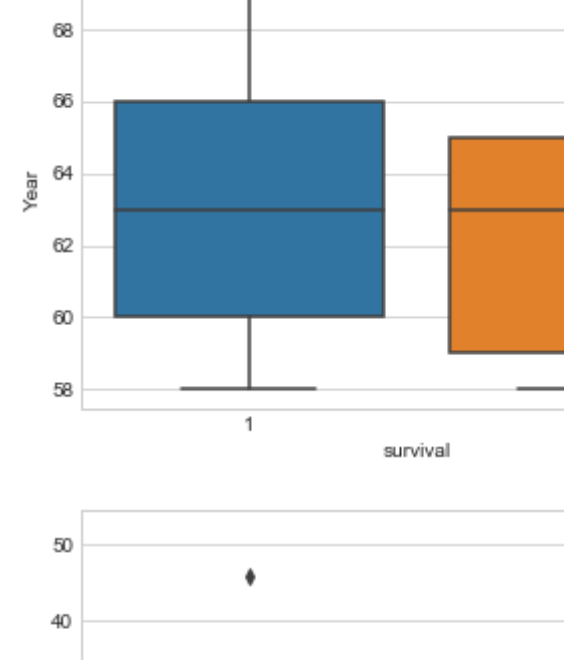
### Observation:-

- 1.CDF of both survival>=5 (i.e survival=1) and survival<5 (i.e survival =2) almost overlap

22.0% of patients who survived 5 years or more were operated at or less than 41 years of age 20% of the patients who didnot survive for 5 years got operated at or less than 45 years of age

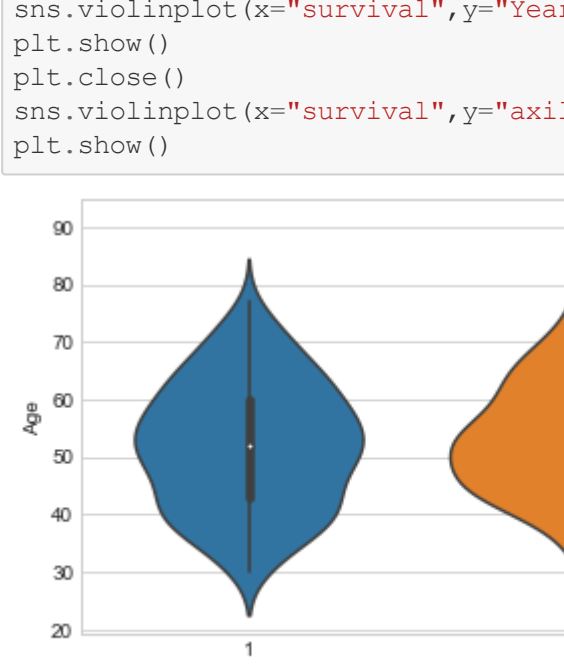
```
In [70]: plt.close()
cnt,bins=np.histogram(data_1["Year"],density=True)
pdf=cnt/(sum(cnt))
cdf=np.cumsum(pdf)
plt.plot(bins[1:],pdf)
plt.plot(bins[1:],cdf)

cnt,bins=np.histogram(data_2["Year"],density=True)
pdf=cnt/(sum(cnt))
cdf=np.cumsum(pdf)
plt.plot(bins[1:],pdf)
plt.plot(bins[1:],cdf)
plt.show()
```

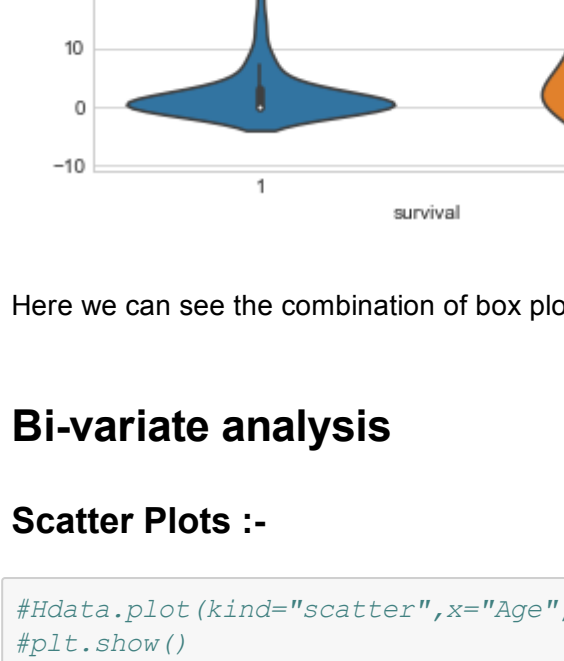


```
In [71]: plt.close()
cnt,bins=np.histogram(data_1["axilnodes"],density=True)
pdf=cnt/(sum(cnt))
cdf=np.cumsum(pdf)
plt.plot(bins[1:],pdf)
plt.plot(bins[1:],cdf)

cnt,bins=np.histogram(data_2["axilnodes"],density=True)
pdf=cnt/(sum(cnt))
cdf=np.cumsum(pdf)
plt.plot(bins[1:],pdf)
plt.plot(bins[1:],cdf)
plt.show()
```



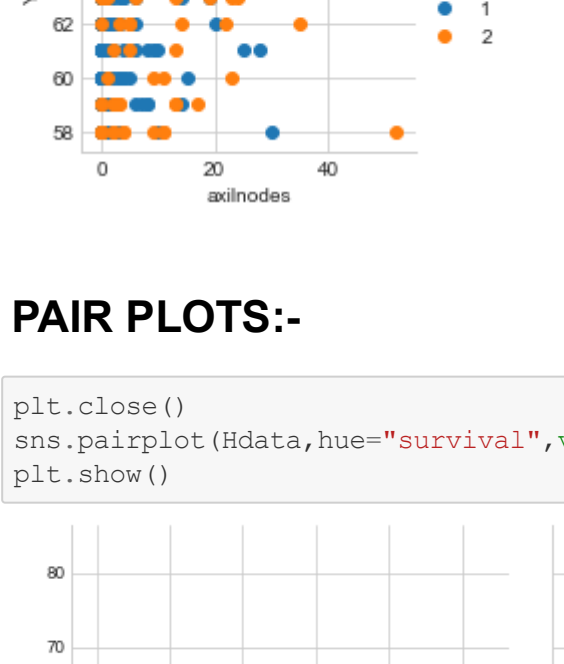
```
In [77]: plt.close()
sns.boxplot(x="survival",y="Age",data=data)
plt.show()
plt.close()
sns.boxplot(x="survival",y="Year",data=data)
plt.show()
plt.close()
sns.boxplot(x="survival",y="axilnodes",data=data)
plt.show()
```



### Observations :-

- 1.Box plots for Age fig(1) are almost similar in terms of ages of 25,50,75 percentile
- 2.Box plots for Year fig(2) are almost similar except 25 % of operations of patients who survived more than 5 years were performed in 1960 or later and 75 % occurred before 1966 whereas 25 % of patients who did not survive 5 years were operated between 1959 and 75% before 1965
- 3.Box plots for auxiliary nodes fig(3) 75% of patients who survived for greater than or equal to 5 years had less than 2 positive auxiliary nodes whereas 75% of patients who did not survive for 5 years had less than 11 positive auxiliary nodes

```
In [78]: plt.close()
sns.violinplot(x="survival",y="Age",data=data)
plt.show()
plt.close()
sns.violinplot(x="survival",y="Year",data=data)
plt.show()
plt.close()
sns.violinplot(x="survival",y="axilnodes",data=data)
plt.show()
```



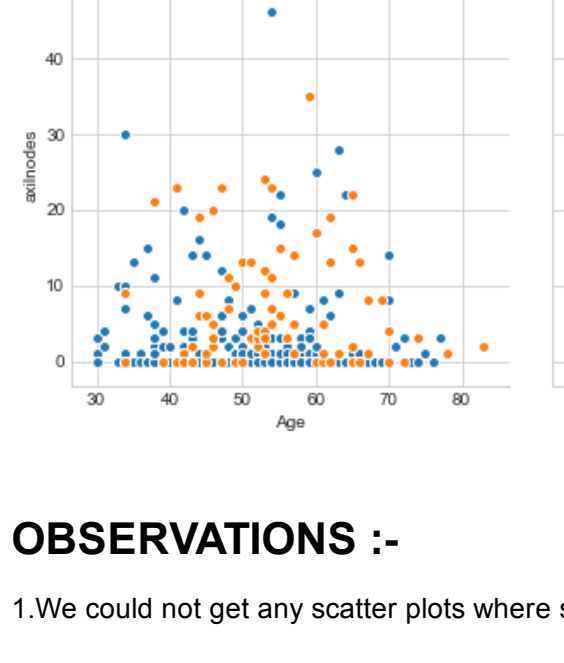
Here we can see the combination of box plot and pdf merged together to form violin plots

### Bi-variate analysis

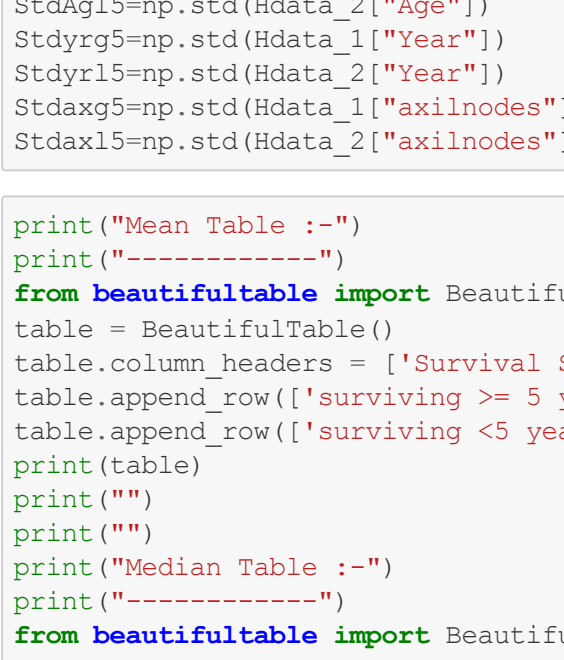
#### Scatter Plots :-

```
In [11]: #data.plot(kind='scatter',x="Age",y="survival")
#plt.show()

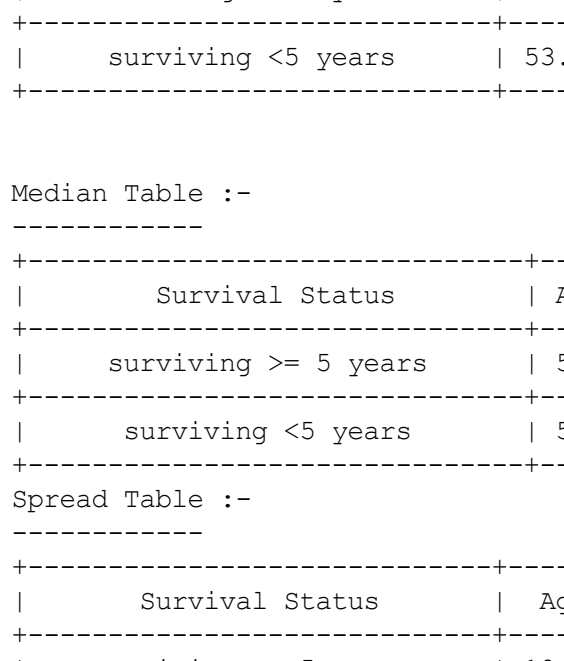
In [22]: sns.set_style("whitegrid");
sns.FacetGrid(data,hue="survival")\
.map(plt.scatter,'Age','Year')\
.add_legend()
plt.show()
```



```
In [6]: sns.set_style("whitegrid");
sns.FacetGrid(data,hue="survival")\
.map(plt.scatter,'Age','Year')\
.add_legend()
plt.show()
```



```
In [8]: sns.set_style("whitegrid");
sns.FacetGrid(data,hue="survival")\
.map(plt.scatter,'axilnodes','Year')\
.add_legend()
plt.show()
```



### PAIR PLOTS:-

```
In [33]: plt.close()
sns.pairplot(data,hue="survival",vars=['Age', 'Year', 'axilnodes'],size=4)
plt.show()
```



### OBSERVATIONS :-

- 1.We could not get any scatter plots where survival status could be linearly seperable
- 2.In case of pair plots the plots of all the features survival staus were highly overlaped and could not be segregated
- 3.We can conclude that even combination of features was not helpful in classification of patients based on survival

```
In [27]: MnAge5=np.mean(data_1["Age"])
MnAge5=np.mean(data_2["Age"])
MnYr5=np.mean(data_1["Year"])
MnYr5=np.mean(data_2["Year"])
Mnax5=np.mean(data_1["axilnodes"])
Mnax5=np.mean(data_2["axilnodes"])
MnAge5=np.median(data_1["Age"])
MnYr5=np.median(data_1["Year"])
Mnax5=np.median(data_1["axilnodes"])
MnAge5=np.median(data_2["Age"])
MnYr5=np.median(data_2["Year"])
Mnax5=np.median(data_2["axilnodes"])

In [38]: StdAge5=np.std(data_1["Age"])
StdAge5=np.std(data_2["Age"])
StdYr5=np.std(data_1["Year"])
StdYr5=np.std(data_2["Year"])
Stdax5=np.std(data_1["axilnodes"])
Stdax5=np.std(data_2["axilnodes"])

In [39]: print("Mean Table :-")
print("-----")
from beautifultable import BeautifulTable
table = BeautifulTable()
table.column_headers = ['Survival Status','Age','Year','No of Auxillary nodes']
table.append_row(['surviving >= 5 years',MnAge5,MnYr5,Mnax5])
table.append_row(['surviving <5 years',MnAge5,MnYr5,Mnax5])
print(table)
print("")
print("Median Table :-")
print("-----")
from beautifultable import BeautifulTable
table = BeautifulTable()
table.column_headers = ['Survival Status','Age','Year','No of Auxillary nodes']
table.append_row(['surviving >= 5 years',MnAge5,MnYr5,Mnax5])
table.append_row(['surviving <5 years',MnAge5,MnYr5,Mnax5])
print(table)
print("Spread Table :-")
print("-----")
table = BeautifulTable()
table.column_headers = ['Survival Status','Age','Year','No of Auxillary nodes']
table.append_row(['surviving >= 5 years',MnAge5,MnYr5,Mnax5])
table.append_row(['surviving <5 years',StdAge5,StdYr5,Stdax5])
print(table)
```

Mean Table :-

	Survival Status	Age	Year	No of Auxillary nodes
1	surviving >= 5 years	52.018	62.862	2.791
2	surviving <5 years	53.079	62.827	7.457

Median Table :-

	Survival Status	Age	Year	No of Auxillary nodes
1	surviving >= 5 years	52.0	63.0	0.0
2	surviving <5 years	53.0	63.0	4.0

Spread Table :-

	Survival Status	Age	Year	No of Auxillary nodes
1	surviving >= 5 years	10.988	3.216	5.857
2	surviving <5 years	10.104	3.321	9.129

### Observations:-

1. Central tendency of Age and Year are quite similar for both the survival categories
2. There is a shift in central tendencies of No of positive aux nodes of both survival categories
3. Although we can see shift in central tendencies for No of Auxiliary nodes but this wont be suitable for classification of survival status as seen from the spread table <5 years has a larger spread and hence also seen by box plots more than 50 % of the data overlap and there is no linear biforcation visible