

3. Plotting for EDA For Haberman Dataset

```
In [1]: 1 import pandas as pd
        2 import seaborn as sns
        3 import matplotlib.pyplot as plt
        4 import numpy as np
        5
        6
        7 haberman = pd.read_csv("haberman.csv")
        8
        9
```

```
In [2]: 1 haberman.head()
```

Out[2]:

	age	year	nodes	status
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

```
In [3]: 1 print (haberman.shape)
        2 #Haberman dataset has 306 rows and 4 columns

(306, 4)
```

```
In [4]: 1 print (haberman.columns)
        2
        3 # Survival status (class attribute) 1 = the patient survived 5 years or longer

Index(['age', 'year', 'nodes', 'status'], dtype='object')
```

```
In [5]: 1 haberman["status"].value_counts()
        2
        3 # We can observe that the haberman is imbalanced dataset as the number of po
```

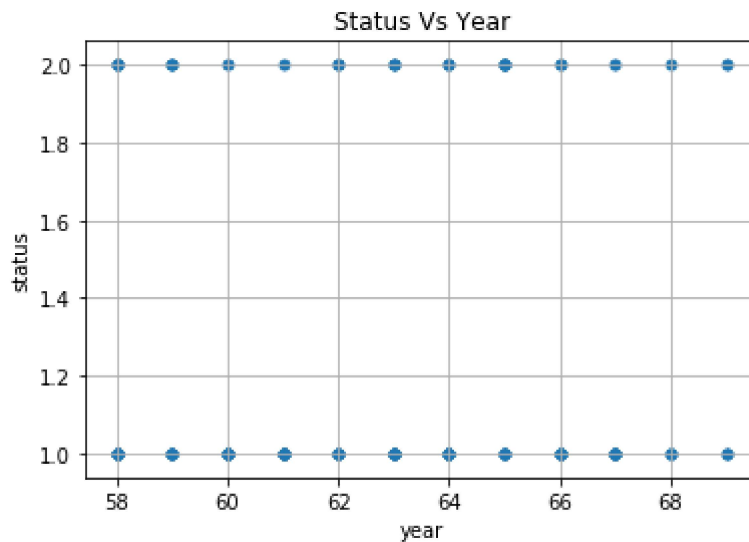
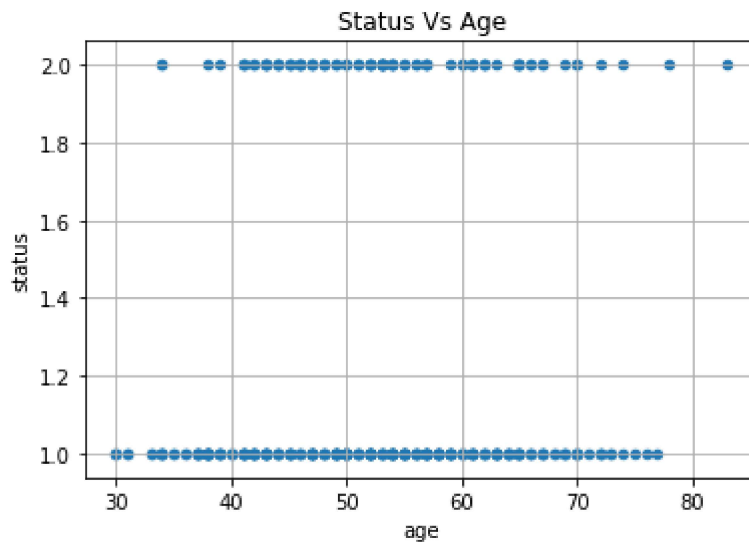
```
Out[5]: 1    225
        2     81
        Name: status, dtype: int64
```

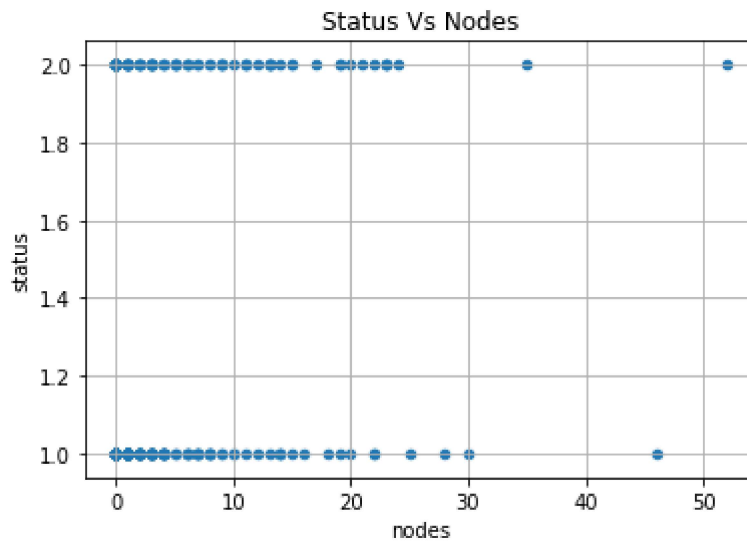
OBJECTIVE

Our objective is to find whether the patient after the treatment for breast cancer survives for more than 5 years or not

2-D Plot

```
In [6]: 1 haberman.plot(kind='scatter', x='age', y='status') ;  
2 plt.title("Status Vs Age")  
3 plt.grid()  
4 plt.show()  
5  
6 haberman.plot(kind='scatter', x='year', y='status') ;  
7 plt.title("Status Vs Year")  
8 plt.grid()  
9 plt.show()  
10  
11 haberman.plot(kind='scatter', x='nodes', y='status') ;  
12 plt.title("Status Vs Nodes")  
13 plt.grid()  
14 plt.show()  
15
```

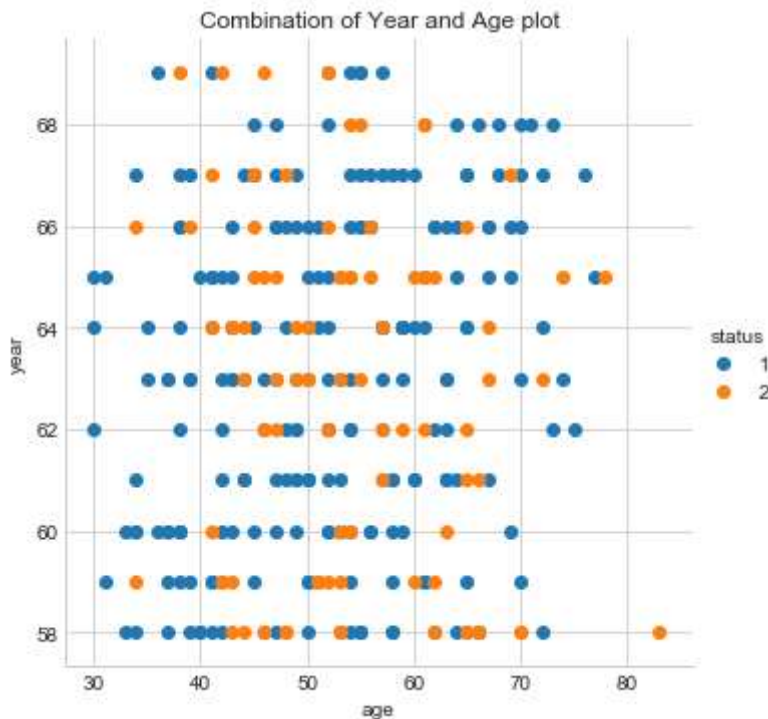




From the graph Age Vs Survival Status, we could observe no kind of dependency on basis of age alone as every year patients died for less than 5 years and survived for more than 5 years. From the graph Year Vs Survival Status, we could observe no kind of dependency on basis of year as every year patients died for less than 5 years and survived for more than 5 years. From the graph Nodes Vs Survival Status, we could observe no kind of dependency on basis of nodes as patients died for less than 5 years and survived for more than 5 years for same Number of positive axillary nodes detected.

So, let's check for the combinations of 2 or more columns

```
In [7]: 1 sns.set_style("whitegrid");  
2 sns.FacetGrid(haberman, hue="status", size=5).map(plt.scatter, 'age', 'year')  
3 plt.title("Combination of Year and Age plot")  
4 plt.show();  
5
```

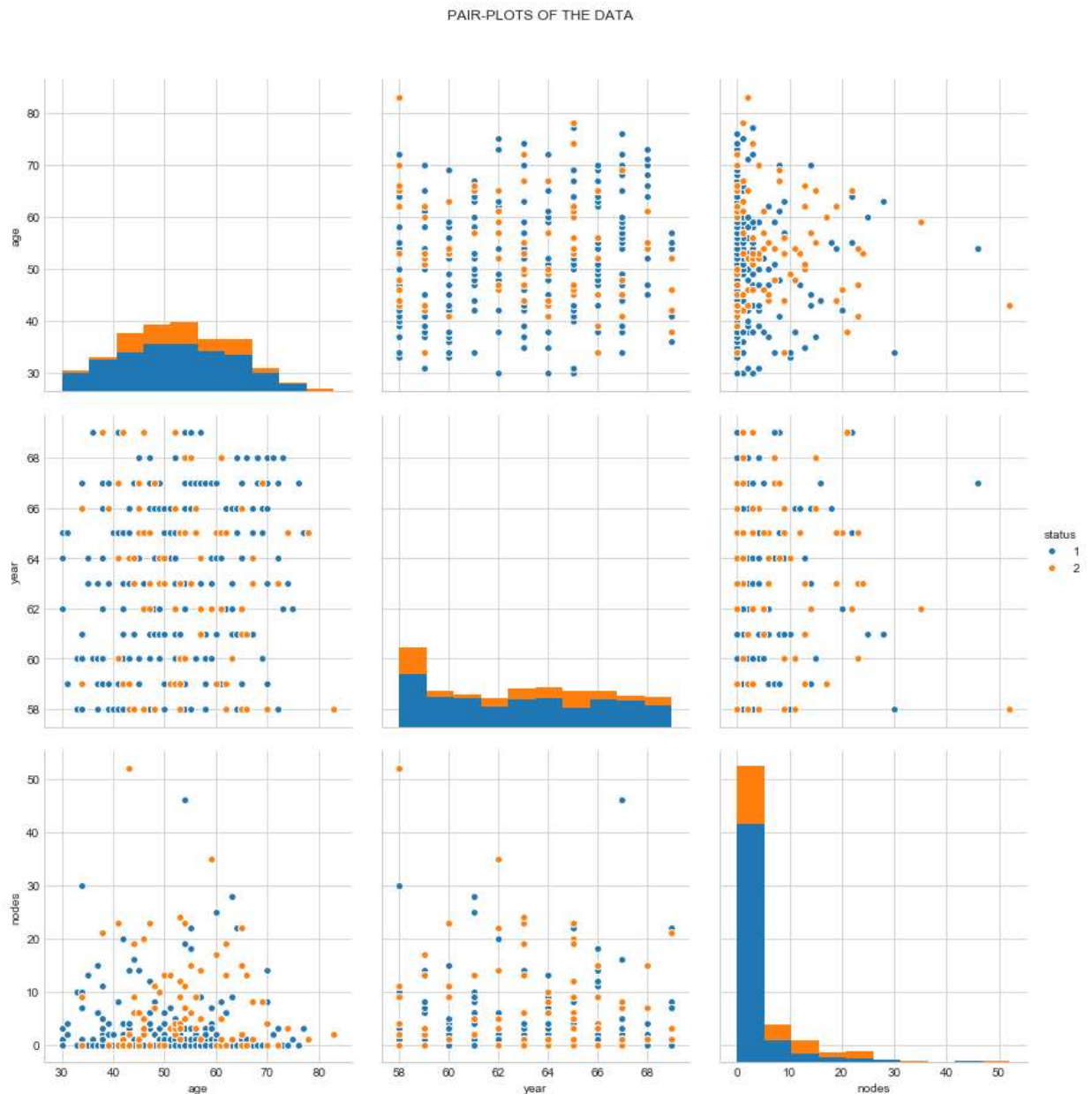


Using the combinations of 2 features, we cant distinguish anything, as mnay of them were overlapping Lets try combining 3 features

```
In [9]: 1 plt.close()
```

Pair-plot

```
In [15]: 1 sns.set_style("whitegrid");
2 g = sns.pairplot(data = haberman, hue="status", size=4, vars = ['age','year',
3 g.fig.suptitle("PAIR-PLOTS OF THE DATA",y = 1.05)
4 plt.show())
```

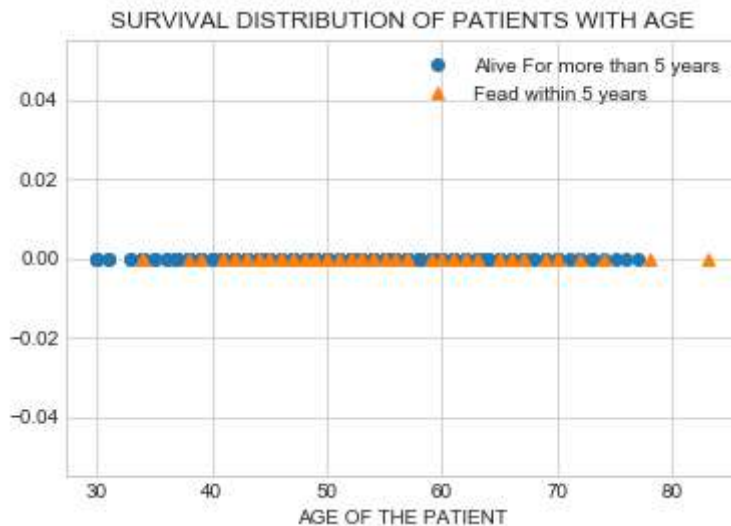


Observations

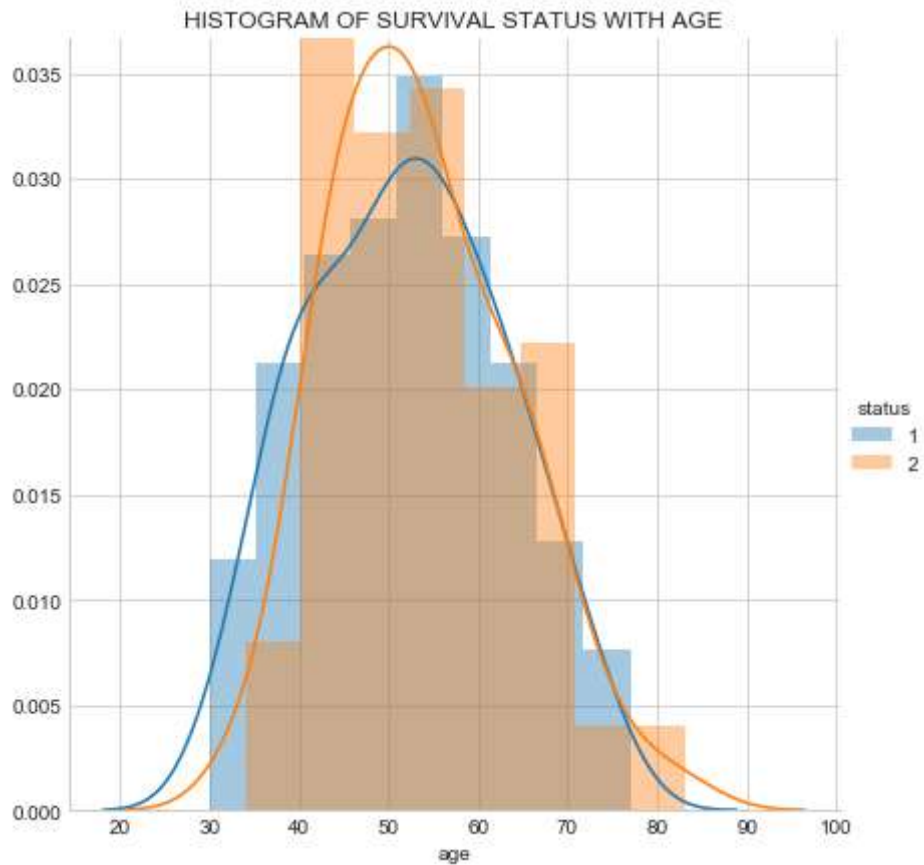
1. age is the most useful features to identify status of a patient.
2. if age ≤ 35 then there is a greater chance of being alive for 5 or more years
3. We can find "lines" and "if-else" conditions to build a simple model to classify the status for the age ≤ 35 .
4. No 2 combinations are that useful. So, we will go with the univariate analysis

(3.4) Histogram, PDF, CDF

```
In [16]: 1 import numpy as np
2 haber_alive = haberman.loc[haberman["status"] == 1];
3 haber_dead = haberman.loc[haberman["status"] == 2];
4 plt.plot(haber_alive["age"], np.zeros_like(haber_alive["age"]), 'o', label = '
5 plt.plot(haber_dead["age"], np.zeros_like(haber_dead["age"]), '^', label = 'Fe
6 plt.title("SURVIVAL DISTRIBUTION OF PATIENTS WITH AGE")
7 plt.xlabel("AGE OF THE PATIENT")
8 plt.legend()
9 plt.show()
```



```
In [18]: 1 import warnings
2 warnings.filterwarnings("ignore")
3 sns.FacetGrid(haberman, hue="status", size=6).map(sns.distplot, "age").add_le
4 plt.title("HISTOGRAM OF SURVIVAL STATUS WITH AGE")
5 plt.show();
6
7
```



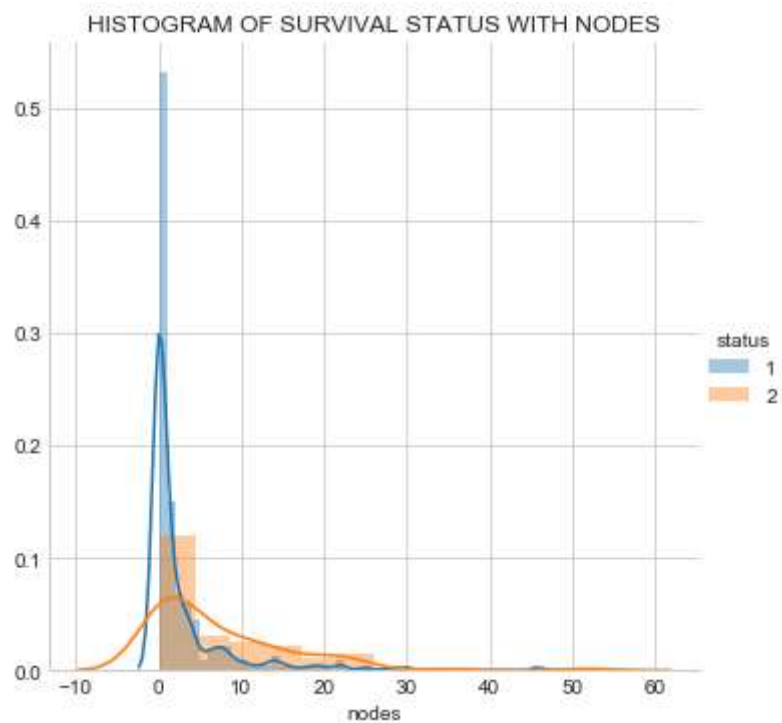
Obersvations

1.If the age is less than the 35, the probability that patient survives more than 5 years is more

2.If the age is more than the 78, the probability that patient dies within 5 years is more

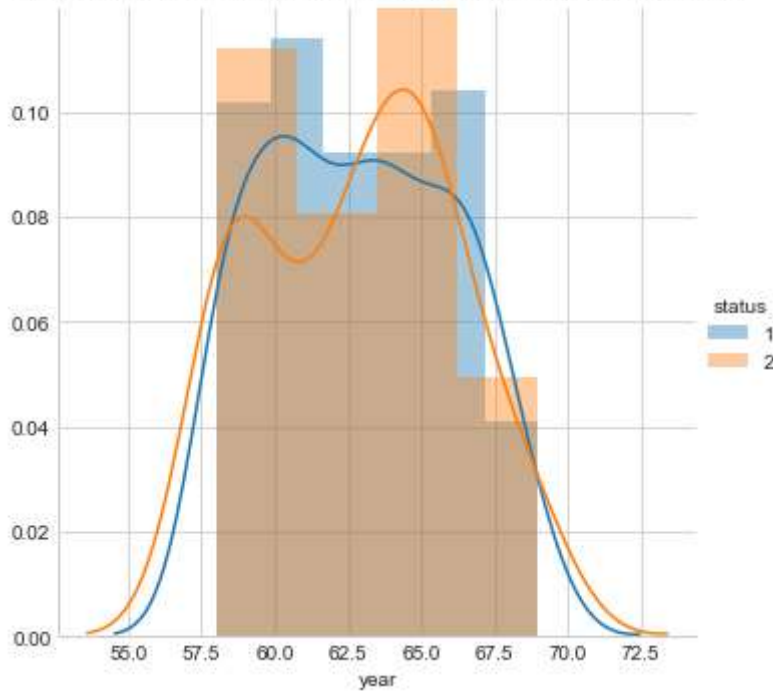
The remaing age gap is unambiguous


```
In [19]: 1 sns.FacetGrid(haberman, hue="status", size=5) \
2         .map(sns.distplot, "nodes") \
3         .add_legend();
4 plt.title("HISTOGRAM OF SURVIVAL STATUS WITH NODES")
5 plt.show();
```



```
In [21]: 1 sns.FacetGrid(haberman, hue="status", size=5) \
2         .map(sns.distplot, "year") \
3         .add_legend();
4 plt.title("HISTOGRAM OF SURVIVAL STATUS WITH YEAR OF OPERATION")
5 plt.show();
```

HISTOGRAM OF SURVIVAL STATUS WITH YEAR OF OPERATION



Obersvations

Nothing can be dervied fro the above 2 histograms. There's so much ambiguity

```

In [24]: 1 counts, bin_edges = np.histogram(haber_alive['age'], bins=10,density = True)
2
3 pdf = counts/(sum(counts))
4 print(pdf);
5 print(bin_edges);
6 cdf = np.cumsum(pdf)
7
8 plt.plot(bin_edges[1:],pdf,label = 'pdf when bin = 10');
9 plt.plot(bin_edges[1:], cdf,label = 'CDF')
10
11 counts, bin_edges = np.histogram(haber_alive['age'], bins=20,density = True)
12
13 pdf = counts/(sum(counts))
14 plt.plot(bin_edges[1:],pdf,label = 'pdf when bin = 20');
15 plt.legend();
16 plt.xlabel("AGE")
17 plt.ylabel("Probability")
18 plt.title("CDF OF SURVIVAL STAUS(ALIVE), PDF WITH 10 BINS AND PDF WITH 20 BIN
19 plt.show();
20

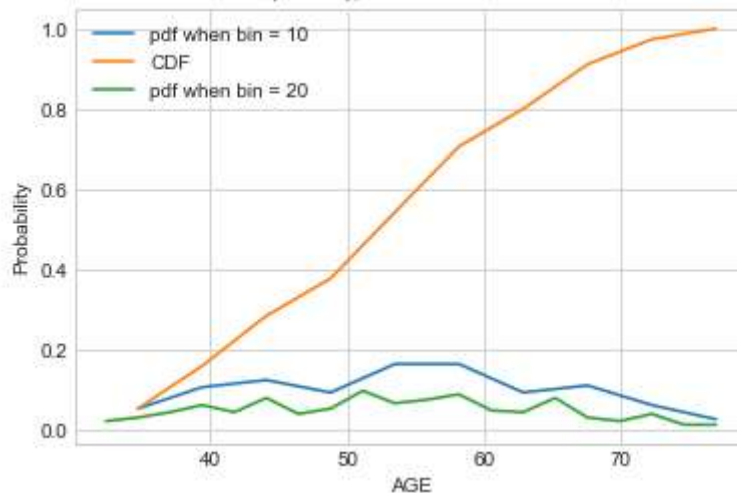
```

```

[0.05333333 0.10666667 0.12444444 0.09333333 0.16444444 0.16444444
0.09333333 0.11111111 0.06222222 0.02666667]
[30.  34.7 39.4 44.1 48.8 53.5 58.2 62.9 67.6 72.3 77. ]

```

CDF OF SURVIVAL STAUS(ALIVE), PDF WITH 10 BINS AND PDF WITH 20 BINS



Obersvations

- 1.If the age is greater than the 55, then there is 60% chance of not surviving more than 5 years
- 2.If the age is greater than the 62, then there is 80% chance of not surviving more than 5 years
- 3.Bin=10 has much better PDF than Bin = 20

```

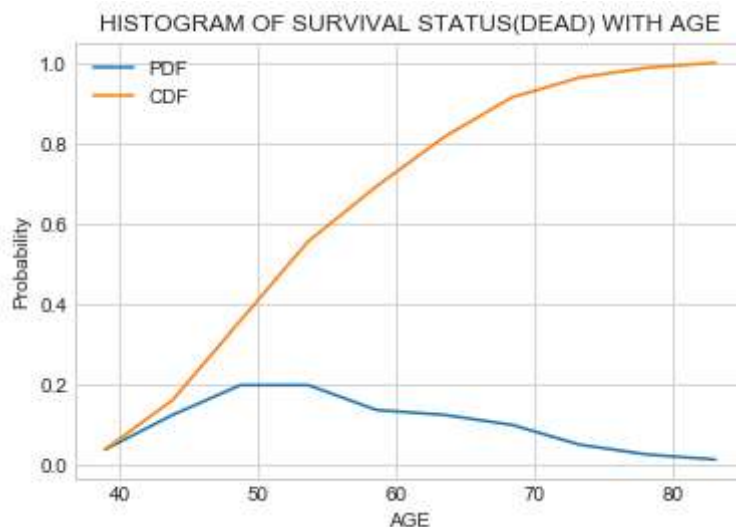
In [23]: 1 counts, bin_edges = np.histogram(haber_dead['age'], bins=10,density = True)
          2
          3 pdf = counts/(sum(counts))
          4 print(pdf);
          5 print(bin_edges)
          6
          7
          8 cdf = np.cumsum(pdf)
          9 plt.plot(bin_edges[1:],pdf,label = "PDF")
         10 plt.plot(bin_edges[1:], cdf,label = "CDF")
         11 plt.legend();
         12 plt.xlabel("AGE")
         13 plt.ylabel("Probability")
         14 plt.title("HISTOGRAM OF SURVIVAL STATUS(DEAD) WITH AGE")
         15 plt.show();

```

```

[0.03703704 0.12345679 0.19753086 0.19753086 0.13580247 0.12345679
 0.09876543 0.04938272 0.02469136 0.01234568]
[34.  38.9 43.8 48.7 53.6 58.5 63.4 68.3 73.2 78.1 83. ]

```



Obersvations

1.The surviving of more than 5 years chances increases till 50 years (20%) and then decreases with growth of age

```

In [25]: 1 counts, bin_edges = np.histogram(haber_alive['age'], bins=10, density = True)
2
3 pdf = counts/(sum(counts))
4 print(pdf);
5 print(bin_edges)
6 cdf = np.cumsum(pdf)
7 plt.plot(bin_edges[1:],pdf,label = "PDF alive")
8 plt.plot(bin_edges[1:], cdf,label = "CDF alive")
9
10
11 # haber_dead
12 counts, bin_edges = np.histogram(haber_dead['age'], bins=10,
13                                density = True)
14 pdf = counts/(sum(counts))
15 print(pdf);
16 print(bin_edges)
17 cdf = np.cumsum(pdf)
18 plt.plot(bin_edges[1:],pdf,label = "PDF dead")
19 plt.plot(bin_edges[1:], cdf,label = "PDF dead")
20 plt.legend();
21 plt.xlabel("AGE")
22 plt.ylabel("Probability")
23 plt.title("PDF's AND CDF's OF PATIENTS WHO ARE ALIVE AND DEAD IN THE SPAN OF
24 plt.show();

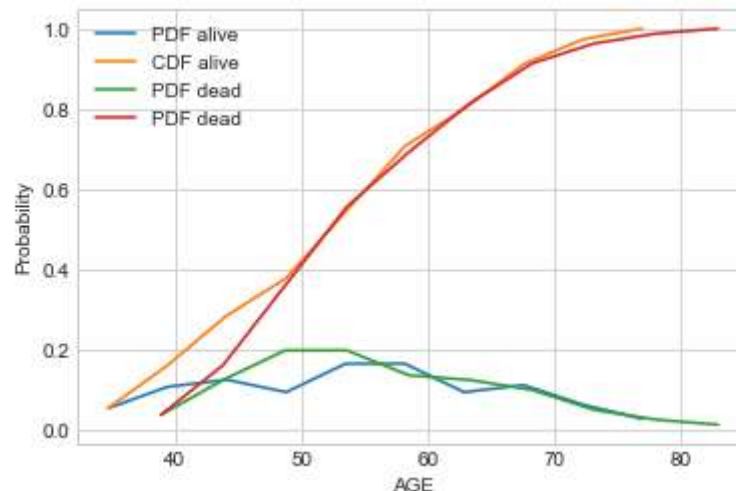
```

```

[0.05333333 0.10666667 0.12444444 0.09333333 0.16444444 0.16444444
 0.09333333 0.11111111 0.06222222 0.02666667]
[30.  34.7 39.4 44.1 48.8 53.5 58.2 62.9 67.6 72.3 77. ]
[0.03703704 0.12345679 0.19753086 0.19753086 0.13580247 0.12345679
 0.09876543 0.04938272 0.02469136 0.01234568]
[34.  38.9 43.8 48.7 53.6 58.5 63.4 68.3 73.2 78.1 83. ]

```

PDF's AND CDF's OF PATIENTS WHO ARE ALIVE AND DEAD IN THE SPAN OF 5 YEARS



Mean, Variance and Std-dev

```
In [16]: 1 #Mean, Variance, Std-deviation,
2 print("Means:")
3 print(np.mean(haber_alive["age"]))
4 #Mean with an outlier.
5 print(np.mean(np.append(haber_alive["age"],5000)));
6 print(np.mean(haber_dead["age"]))
7
8 print("\nStd-dev:");
9 print(np.std(haber_alive["age"]))
10 print(np.std(np.append(haber_alive["age"],5000)))
11 print(np.std(haber_dead["age"]))
```

Means:

52.01777777777778

73.91150442477876

53.67901234567901

Std-dev:

10.98765547510051

328.58884542338734

10.10418219303131

Observation

Mean, variance is changing rapidly for one value

Median, Percentile, Quantile, IQR, MAD

```
In [17]: 1 #Median, Quantiles, Percentiles, IQR.
2 print("\nMedians:")
3 print(np.median(haber_alive["age"]))
4 #Median with an outlier
5 print(np.median(np.append(haber_alive["age"],50)));
6 print(np.median(haber_dead["age"]))
7
8
9 print("\nQuantiles:")
10 print(np.percentile(haber_alive["age"],np.arange(0, 100, 25)))
11 print(np.percentile(haber_dead["age"],np.arange(0, 100, 25)))
12
13 print("\n90th Percentiles:")
14 print(np.percentile(haber_alive["age"],90))
15 print(np.percentile(haber_dead["age"],90))
16
17 from statsmodels import robust
18 print ("\nMedian Absolute Deviation")
19 print(robust.mad(haber_alive["age"]))
20 print(robust.mad(haber_dead["age"]))
21
```

Medians:

52.0

52.0

53.0

Quantiles:

[30. 43. 52. 60.]

[34. 46. 53. 61.]

90th Percentiles:

67.0

67.0

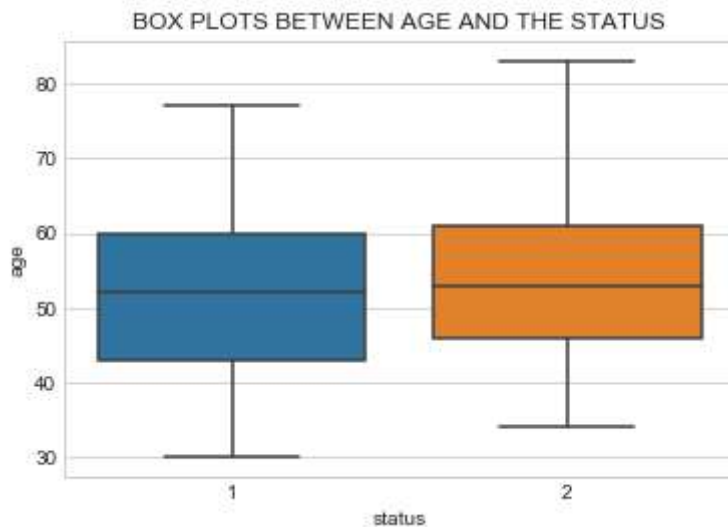
Median Absolute Deviation

13.343419966550417

11.860817748044816

Box plot and Whiskers

```
In [28]: 1 sns.boxplot(x='status',y='age', data=haberman)
2 plt.title("BOX PLOTS BETWEEN AGE AND THE STATUS")
3 plt.show()
```



Observations

Patients who survived for 5 or more has(all values are approx)

1. 50th percentile : 53
2. 25th percentile : 43
3. 75th percentile : 60
4. Whiskers are min max of age of the people of this category

-> IQR = 17

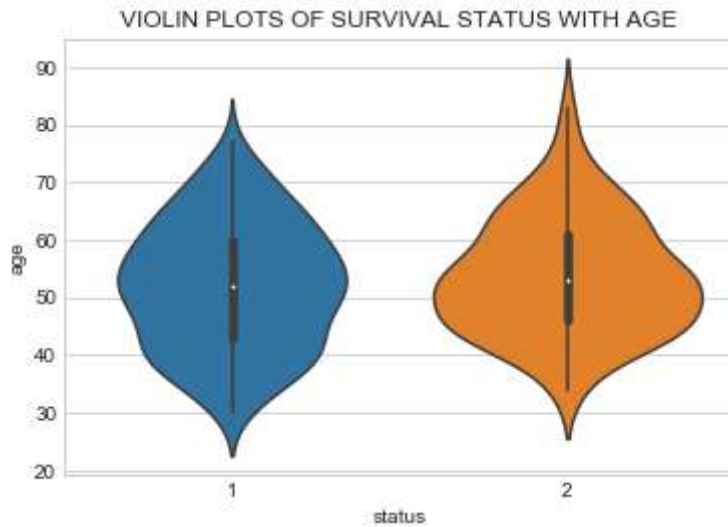
Patients who didn't survived for 5 or more has(all values are approx)

1. 50th percentile : 54
2. 25th percentile : 46
3. 75th percentile : 62
4. Whiskers are min max of age of the people of this category

-> IQR = 16

(3.8) Violin plots

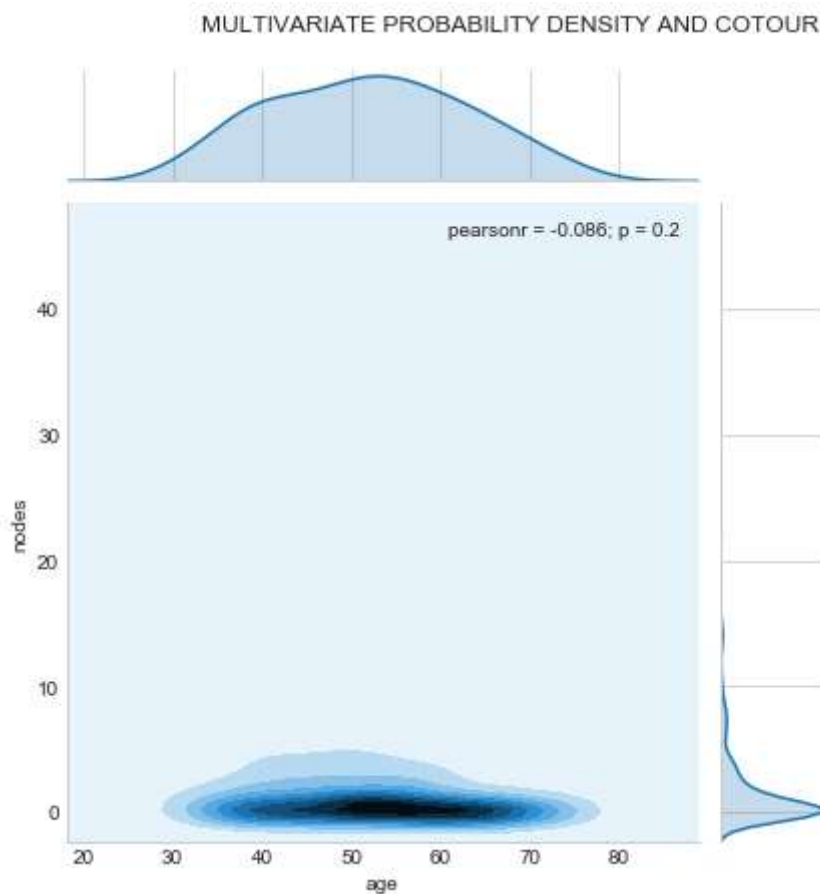

```
In [29]: 1 sns.violinplot(x="status", y="age", data=haberman, size=8)
2 plt.title("VIOLIN PLOTS OF SURVIVAL STATUS WITH AGE")
3 plt.show()
```



1. In the orange plot we can observe that its wider at the age of 50. So we could say that there are larger number of people at the age of 50 would have survived more than 5 years

Multivariate probability density, contour plot.

```
In [38]: 1 #2D Density plot, contours-plot
2 sns.jointplot(x="age", y="nodes", data=haber_alive, kind="kde");
3 plt.title("MULTIVARIATE PROBABILITY DENSITY AND COTOUR PLOT",x = -1.6,y = 1.2)
4 plt.show();
5
```



CONCLUSIONS

1. The Patients survived for more than 5 years are a bit greater in number for patients treated a few years back than those of patients treated many year back So we could say that the technology advancement has helped the patients suffering from the breast cancer
2. Age factor : The patients who were old didn't make for more than 5 years mostly but those who were young survived for more than 5 years