

```
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
import numpy as np

In [3]: haberman = pd.read_csv("haberman (2).csv")

In [3]: haberman

Out[3]:
```

	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
...
301	75	62	1	1
302	76	67	0	1
303	77	65	3	1
304	78	65	1	2
305	83	58	2	2

306 rows × 4 columns

```
In [ ]: #this dataset consist of total list of 306 patients that classified on basis of their age,year of operation ,nodes
#status has numerical value of either 1 (that represents patients survive 5 years or more after operation) and
```

```
In [5]: haberman["status"].value_counts()
```

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

```
In [ ]: # so in dataset we have total of 225 number of patients with survival years or more and 81 no of patients with
```

```
In [6]: habex = haberman["status"]=="1"
        habex
```


```
Out[6]: 0      False
        1      False
        2      False
        3      False
        4      False
        ...
        301    False
        302    False
        303    False
        304    False
        305    False
        Name: status, Length: 306, dtype: bool
```

```
In [ ]: #As status 1 represents patients surviving 5 years or more and status 2 represents patients surviving 5 years or less
# So out of 306 patients we can say 225 patients survived 5 years or more and 81 patients survived 5 years or less
```

```
In [ ]:
```

```
In [5]: sns.FacetGrid(haberman, hue="status", size=4) \
        .map(plt.scatter, "age", "year") \
        .add_legend();
        plt.show();
```

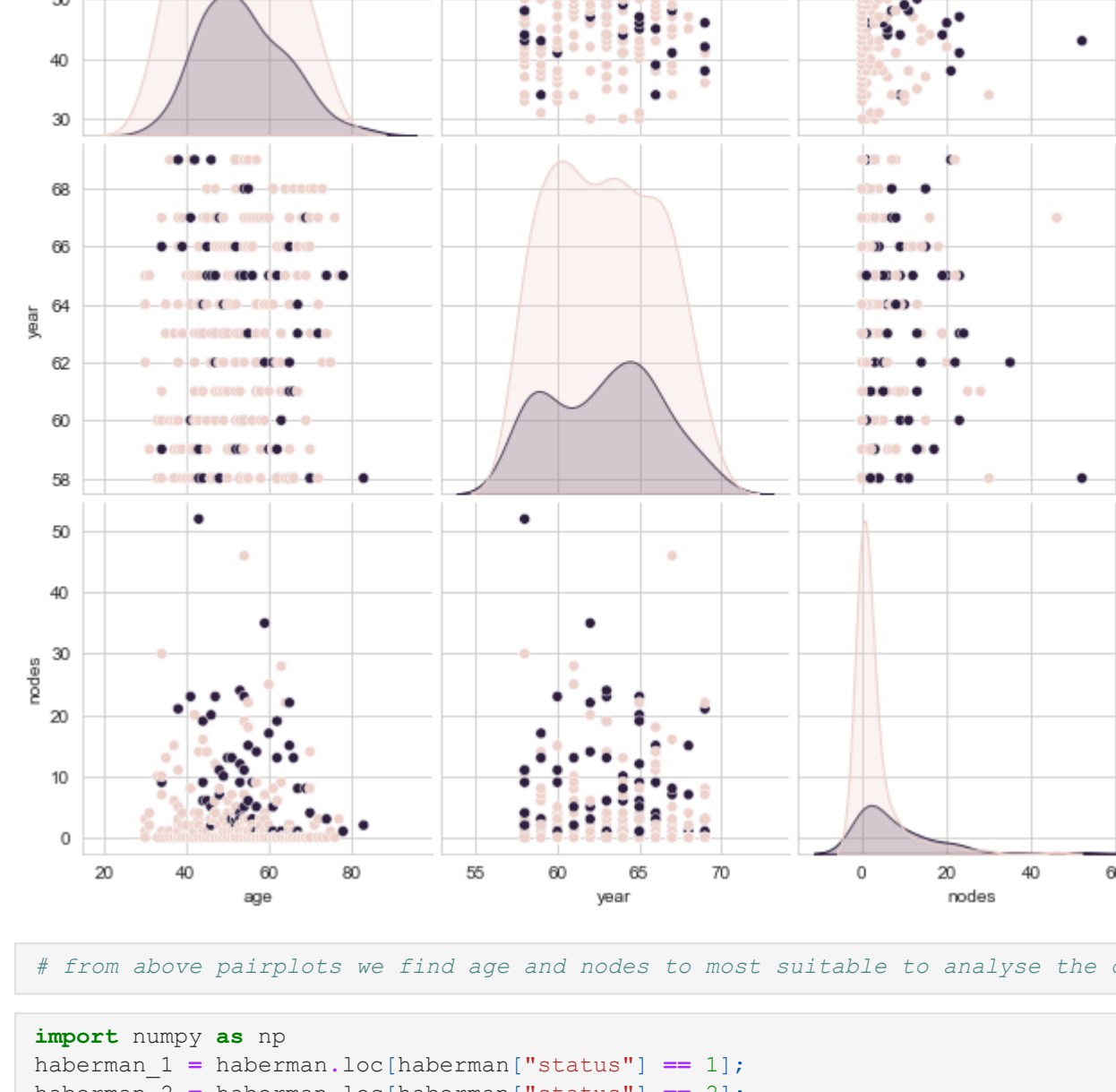
/Users/apple/opt/anaconda3/lib/python3.8/site-packages/seaborn/axisgrid.py:316: UserWarning: The 'size' parameter has been renamed to 'height'; please update your code.
warnings.warn(msg, UserWarning)



```
In [ ]:
```

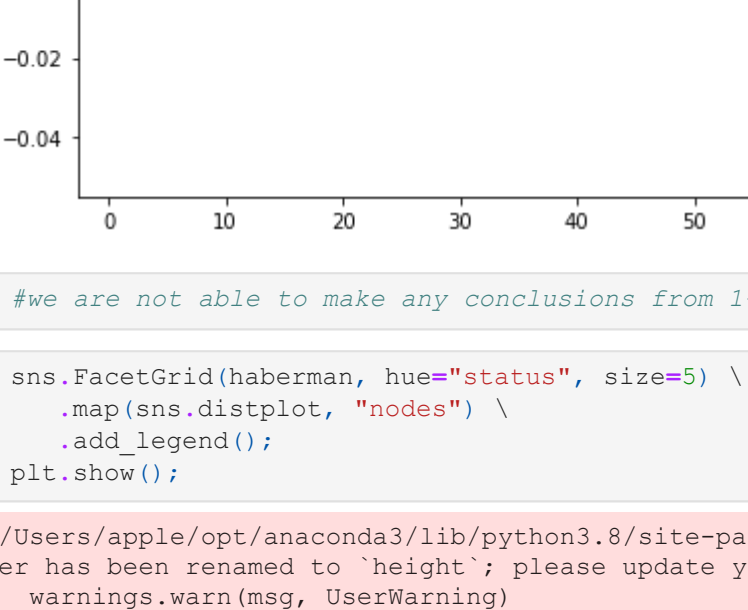
```
In [7]: sns.set_style("whitegrid");
        sns.pairplot(haberman, hue="status", size=3);
        plt.show();
```

/Users/apple/opt/anaconda3/lib/python3.8/site-packages/seaborn/axisgrid.py:1912: UserWarning: The 'size' parameter has been renamed to 'height'; please update your code.
warnings.warn(msg, UserWarning)



```
In [ ]: # From above pairplots we find age and nodes to most suitable to analyse the data further as it has least overlap
```

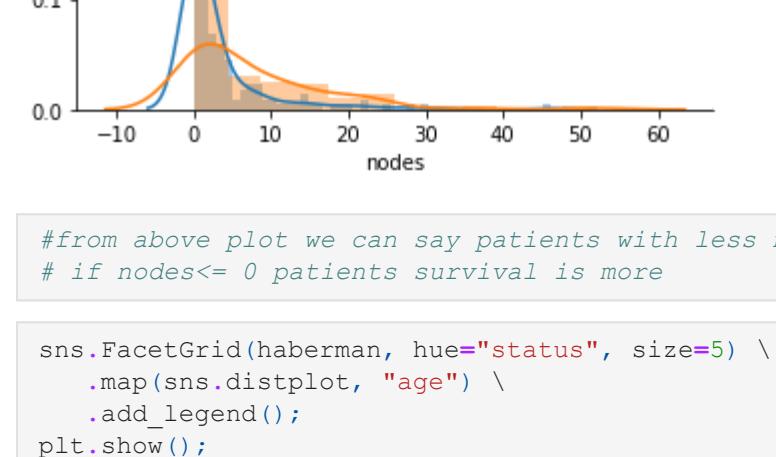
```
In [7]: import numpy as np
        haberman_1 = haberman.loc(haberman["status"] == 1);
        haberman_2 = haberman.loc(haberman["status"] == 2);
        plt.plot(haberman_1["nodes"], np.zeros_like(haberman_1["nodes"]), 'o')
        plt.plot(haberman_2["nodes"], np.zeros_like(haberman_2["nodes"]), 'o')
        #plt.plot(iris_virginica["petal_length"], np.zeros_like(iris_virginica["petal_length"]), 'o')
        plt.show();
```



```
In [ ]: #we are not able to make any conclusions from 1-D scatter plot as most of points are overlapping each other
```

```
In [8]: sns.FacetGrid(haberman, hue="status", size=5) \
        .map(sns.distplot, "nodes") \
        .add_legend();
        plt.show();
```

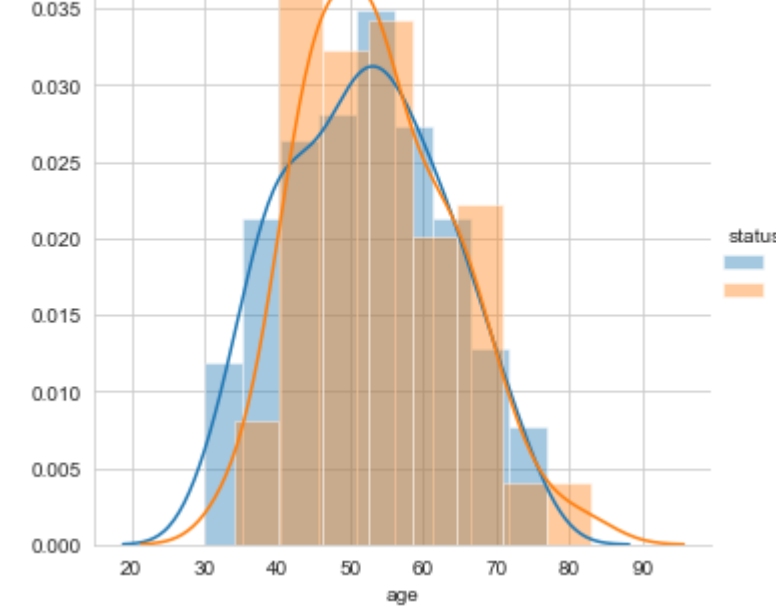
/Users/apple/opt/anaconda3/lib/python3.8/site-packages/seaborn/axisgrid.py:316: UserWarning: The 'size' parameter has been renamed to 'height'; please update your code.
warnings.warn(msg, UserWarning)
/Users/apple/opt/anaconda3/lib/python3.8/site-packages/seaborn/distributions.py:2551: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)
/Users/apple/opt/anaconda3/lib/python3.8/site-packages/seaborn/distributions.py:2551: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)



```
In [ ]: #from above plot we can say patients with less number of nodes tend to survive more
# if nodes<= 0 patients survival is more
```

```
In [15]: sns.FacetGrid(haberman, hue="status", size=5) \
        .map(sns.distplot, "age") \
        .add_legend();
        plt.show();
```

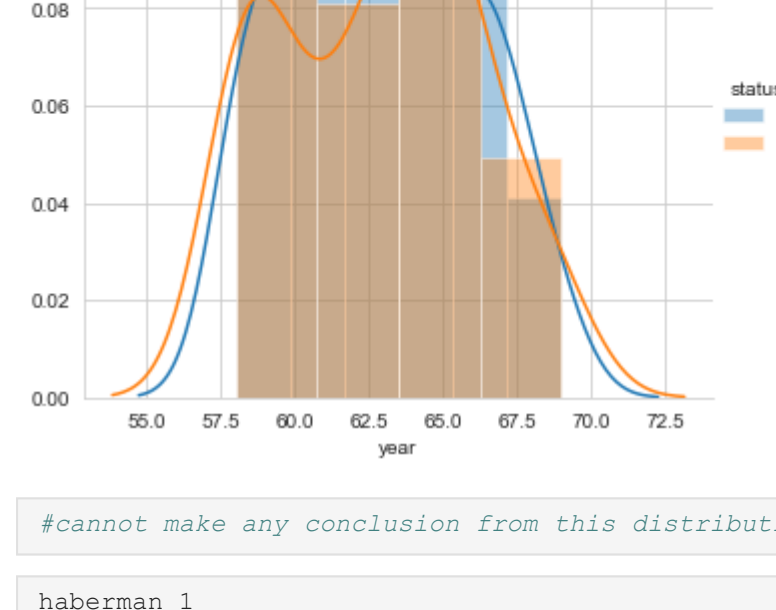
/Users/apple/opt/anaconda3/lib/python3.8/site-packages/seaborn/axisgrid.py:316: UserWarning: The 'size' parameter has been renamed to 'height'; please update your code.
warnings.warn(msg, UserWarning)
/Users/apple/opt/anaconda3/lib/python3.8/site-packages/seaborn/distributions.py:2551: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)
/Users/apple/opt/anaconda3/lib/python3.8/site-packages/seaborn/distributions.py:2551: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)



```
In [ ]: # From this plot though there is no clear saturation of age still we can make a rough prediction as
# if age<35 patients usually survive more than 5 years
# else if age > 75 patients usually survive less than 5 years
```

```
In [16]: sns.FacetGrid(haberman, hue="status", size=5) \
        .map(sns.distplot, "year") \
        .add_legend();
        plt.show();
```

/Users/apple/opt/anaconda3/lib/python3.8/site-packages/seaborn/axisgrid.py:316: UserWarning: The 'size' parameter has been renamed to 'height'; please update your code.
warnings.warn(msg, UserWarning)
/Users/apple/opt/anaconda3/lib/python3.8/site-packages/seaborn/distributions.py:2551: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)
/Users/apple/opt/anaconda3/lib/python3.8/site-packages/seaborn/distributions.py:2551: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)



```
In [13]: #cannot make any conclusion from this distribution plot as it is overlapping at every point
```

```
In [9]: haberman_1
```

```
Out[9]:
```

	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
...
298	73	68	0	1
300	74	63	0	1
301	75	62	1	1
302	76	67	0	1
303	77	65	3	1

225 rows × 4 columns

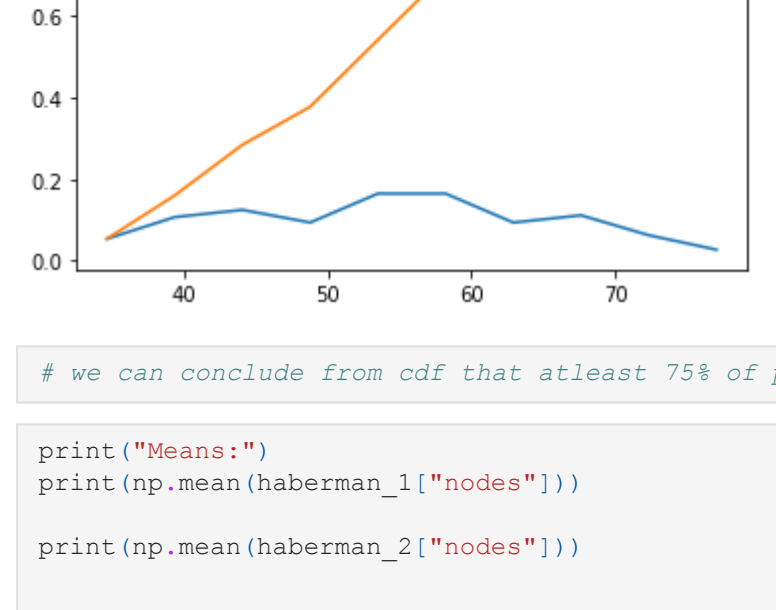
```
In [10]: counts, bin_edges = np.histogram(haberman_1["age"], bins=10,
        density = True);
        print(counts)
        pdf = counts/(sum(counts))

        print(pdf);
        print(bin_edges)

        #compute CDF
        cdf = np.cumsum(pdf)
        plt.plot(bin_edges[1:],pdf)
        plt.plot(bin_edges[1:],cdf)

        plt.show();
```

[0.01134752 0.02269504 0.02647754 0.01985816 0.03498818 0.03498818
0.01985816 0.02364066 0.01323877 0.00567376]
[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.]



```
In [ ]: # we can conclude from cdf that atleast 75% of patients with age <=60 tend to live more than five years
```

```
In [29]: print("Means:")
        print(np.mean(haberman_1["nodes"]))

        print(np.mean(haberman_2["nodes"]))

        print("\nStd-dev:")
        print(np.std(haberman_1["nodes"]))
        print(np.std(haberman_2["nodes"]))
```

Means:
2.7911111111111113
7.45679012345679
Std-dev:
5.857258449412131
9.128776076761632

```
In [ ]:
```

```
In [11]: print("\nMedians:")
        print(np.median(haberman_1["nodes"]))
        #Median with an outlier
        print(np.median(np.append(haberman_1["nodes"],50)));
        print(np.median(haberman_2["nodes"]))
```

Medians:
0.0
4.0
4.0

```
In [ ]: # average node size of patients with survival more than 5 years is 0 .
# average node size of patients with survival less than 5 years is 4.
```

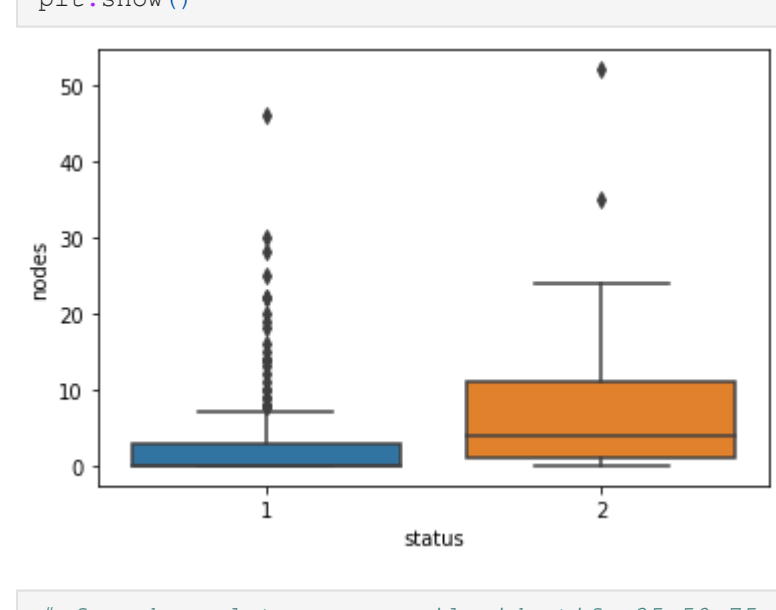
```
In [12]: print("\nQuantiles:")
        print(np.percentile(haberman_1["nodes"],np.arange(0, 100, 25)))
        print(np.percentile(haberman_2["nodes"],np.arange(0, 100, 25)))

        print("\n90th Percentiles:")
        print(np.percentile(haberman_1["nodes"],90))
        print(np.percentile(haberman_2["nodes"],90))
```

Quantiles:
[0. 0. 0. 3.]
[0. 1. 4. 11.]
90th Percentiles:
8.0
20.0

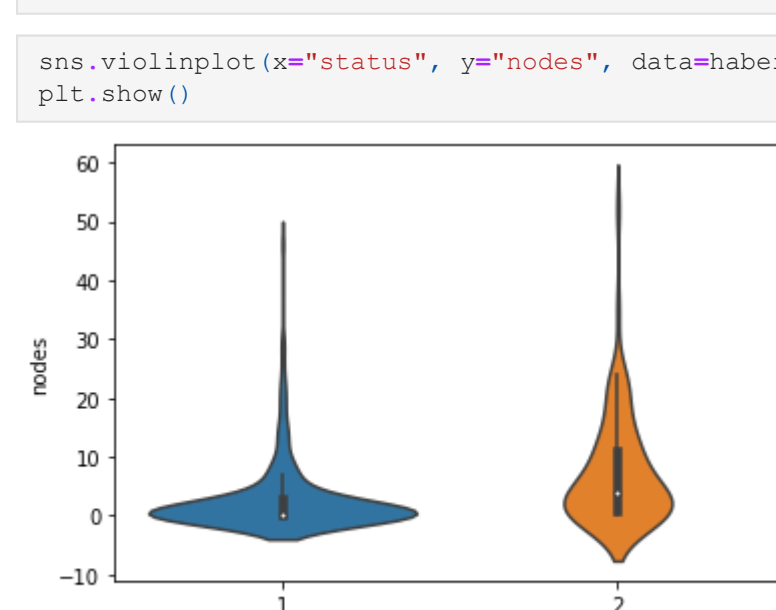
```
In [ ]: # with quantiles we can say 50 % of patients having survival more than 5 years have node 0 .
# atleast 75% of patients having survival more than 5 years have node less than 3
# 90 % of patients with survival more than 5 years have node <= 8
# with quantiles we can also say 50 % of patients with survival less than 5 years have node <=4 and 75% of patients with survival less than 5 years have node <= 20
```

```
In [13]: sns.boxplot(x="status",y="nodes", data=haberman)
        plt.show();
```



```
In [ ]: # from box plot we can easily identify 25-50-75 percentile of values .
# 75 % of patients with survival more than 5 years have node size <= 2
# only 25 % of patients with survival less than 5 years have node size <=1
and 75 % of patients with survival less than 5 years have node size <=11
# 50 % of patients with survival less than 5 years have node size <=5
```

```
In [14]: sns.violinplot(x="status", y="nodes", data=haberman, size=8)
        plt.show();
```



```
In [ ]: # violin plot are combination of histogram , pdf and boxplot .
# we can say patients with survival more than 5 years have node range 0<=node <=9.
# we can say patients with survival less than 5 years have node range 0<=node <=25.
```