

```
In [94]: import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
import pandas as pd
```

```
In [95]: hbrmn = pd.read_csv("haberman.csv") # loading data
print(hbrmn.shape) # print the shape to understand the row and columns

(306, 4)
```

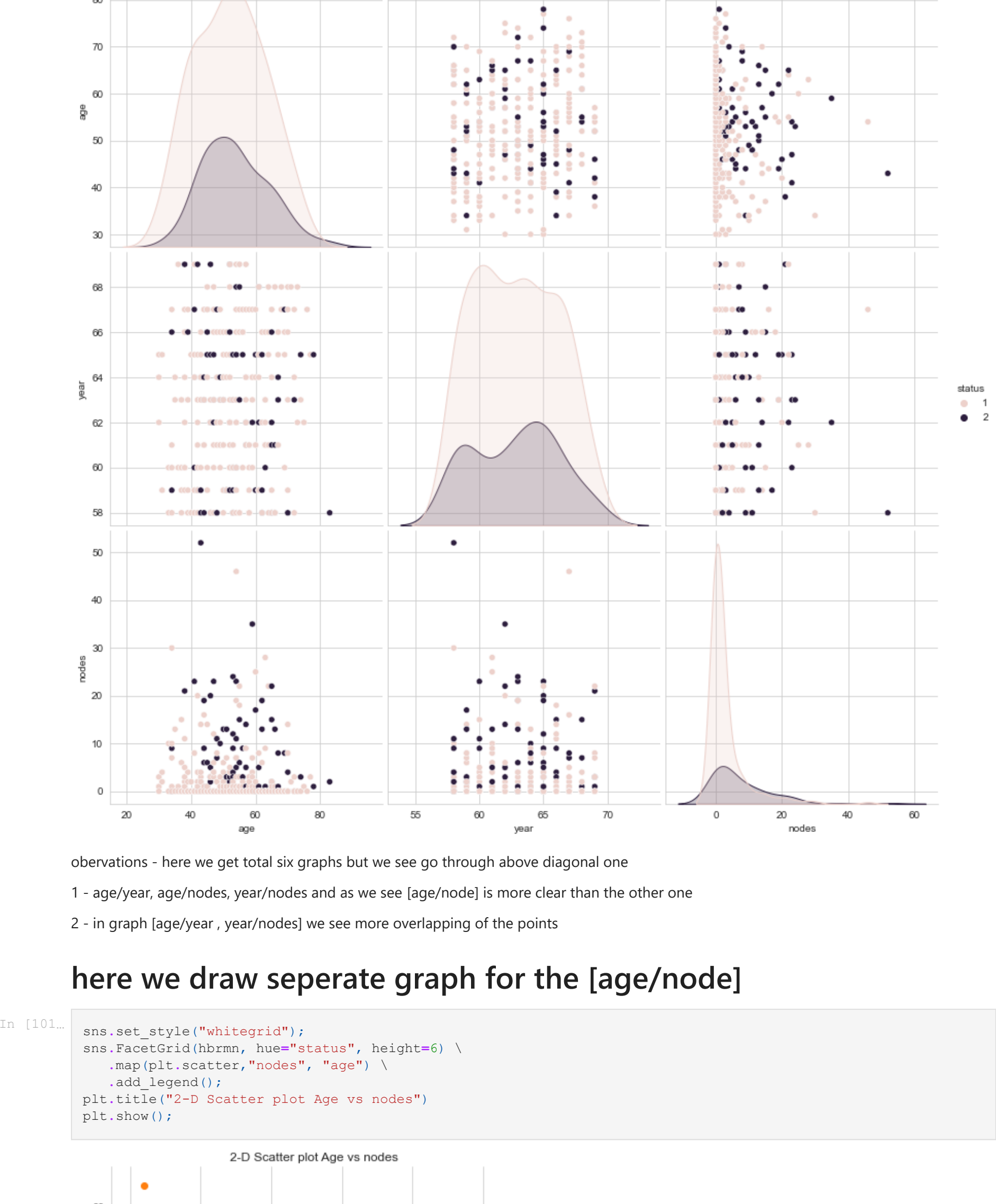
```
In [98]: print(hbrmn.columns) # print the all features

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

```
In [99]: hbrmn["status"].value_counts() # count the no data in each class

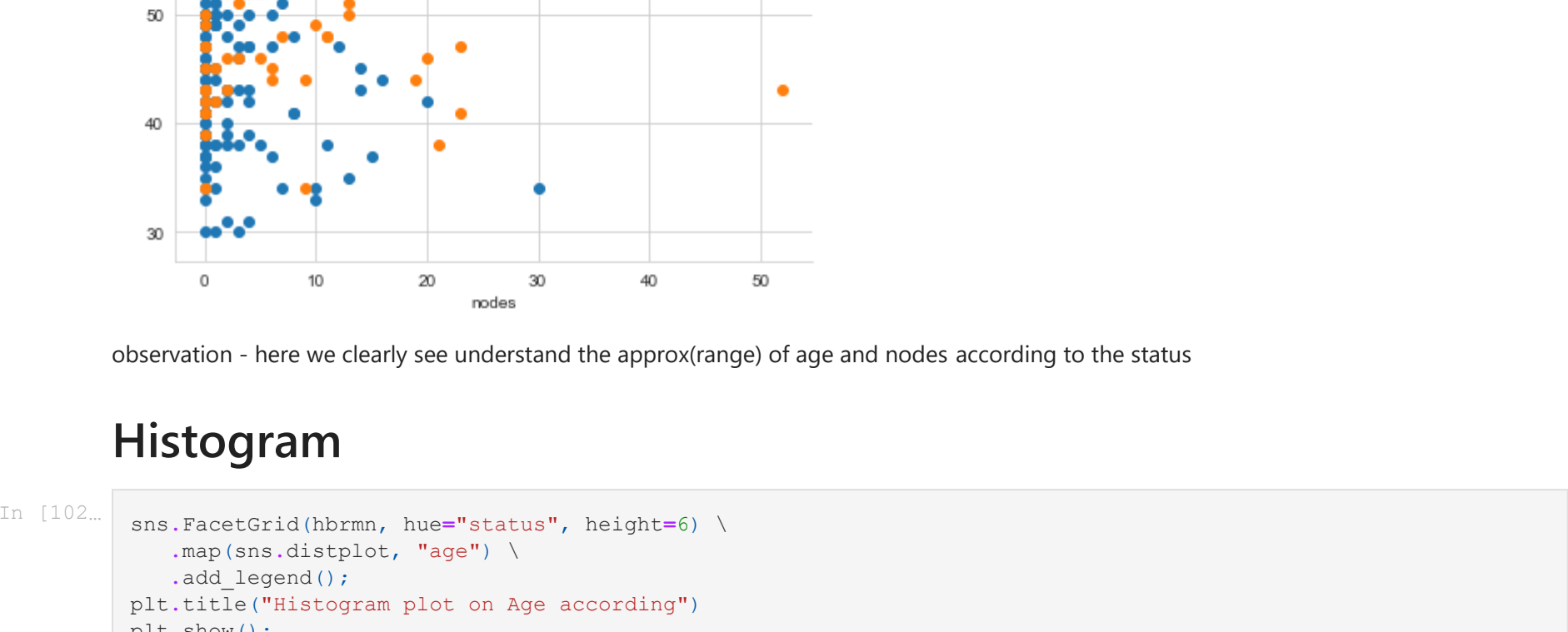
Out[99]: 1    225
         2     81
         Name: status, dtype: int64
```

plotting the pairplot, it give graphical representation of each feature with respect to other one



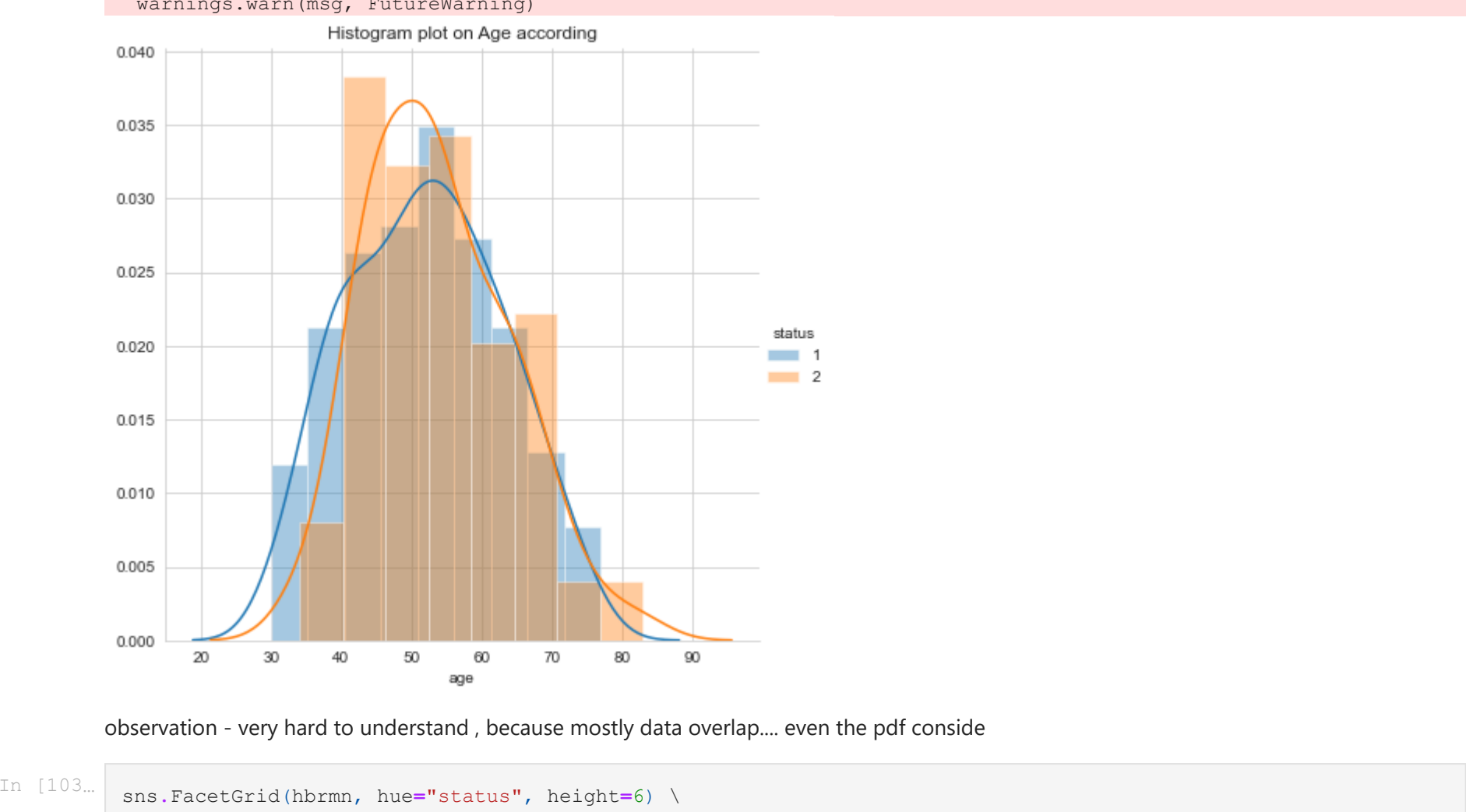
observations - here we get total six graphs but we see go through above diagonal one
1 - age/year, age/nodes, year/nodes and as we see [age/node] is more clear than the other one
2 - in graph [age/year, year/nodes] we see more overlapping of the points

here we draw separate graph for the [age/node]

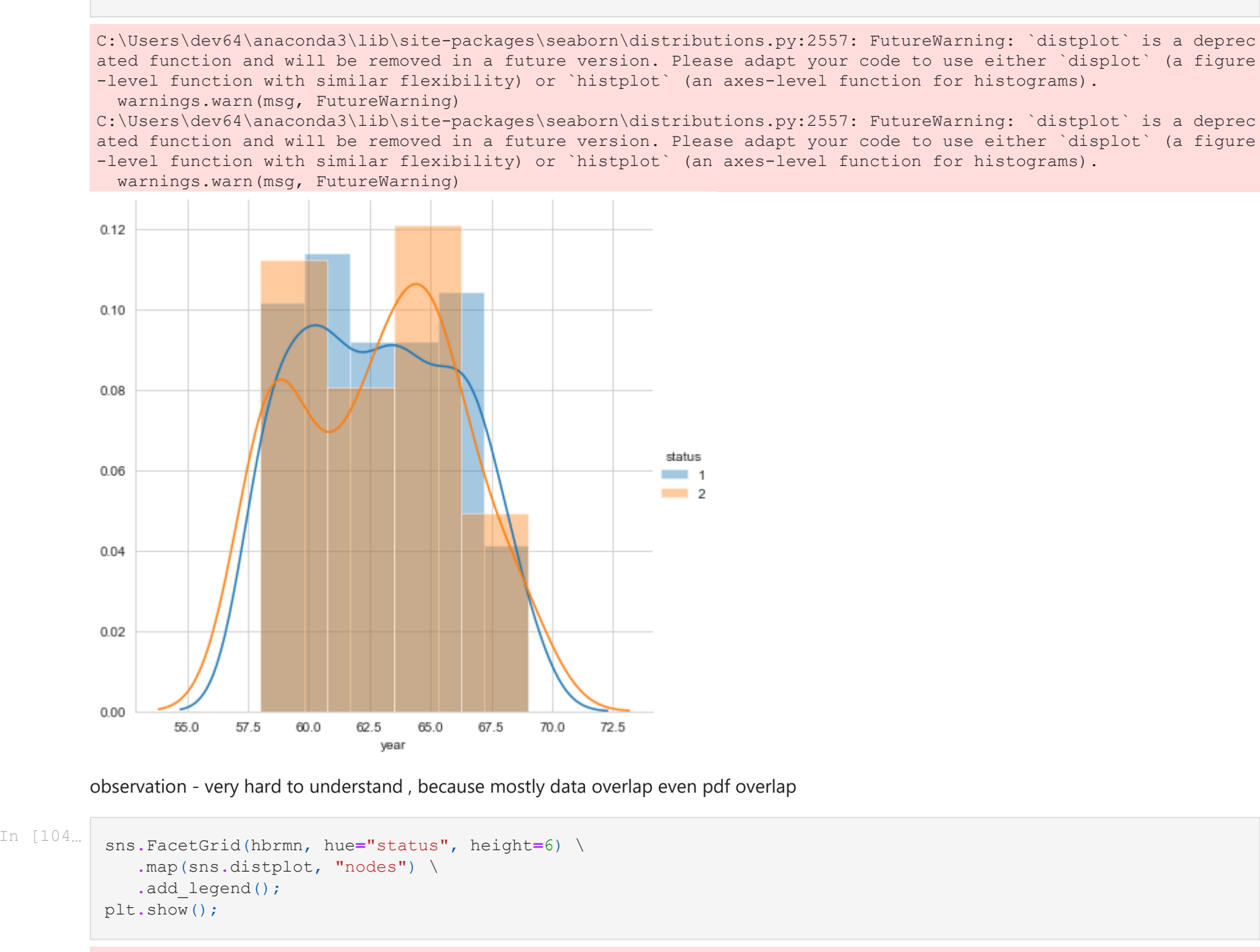


observation - here we clearly see understand the approx(range) of age and nodes according to the status

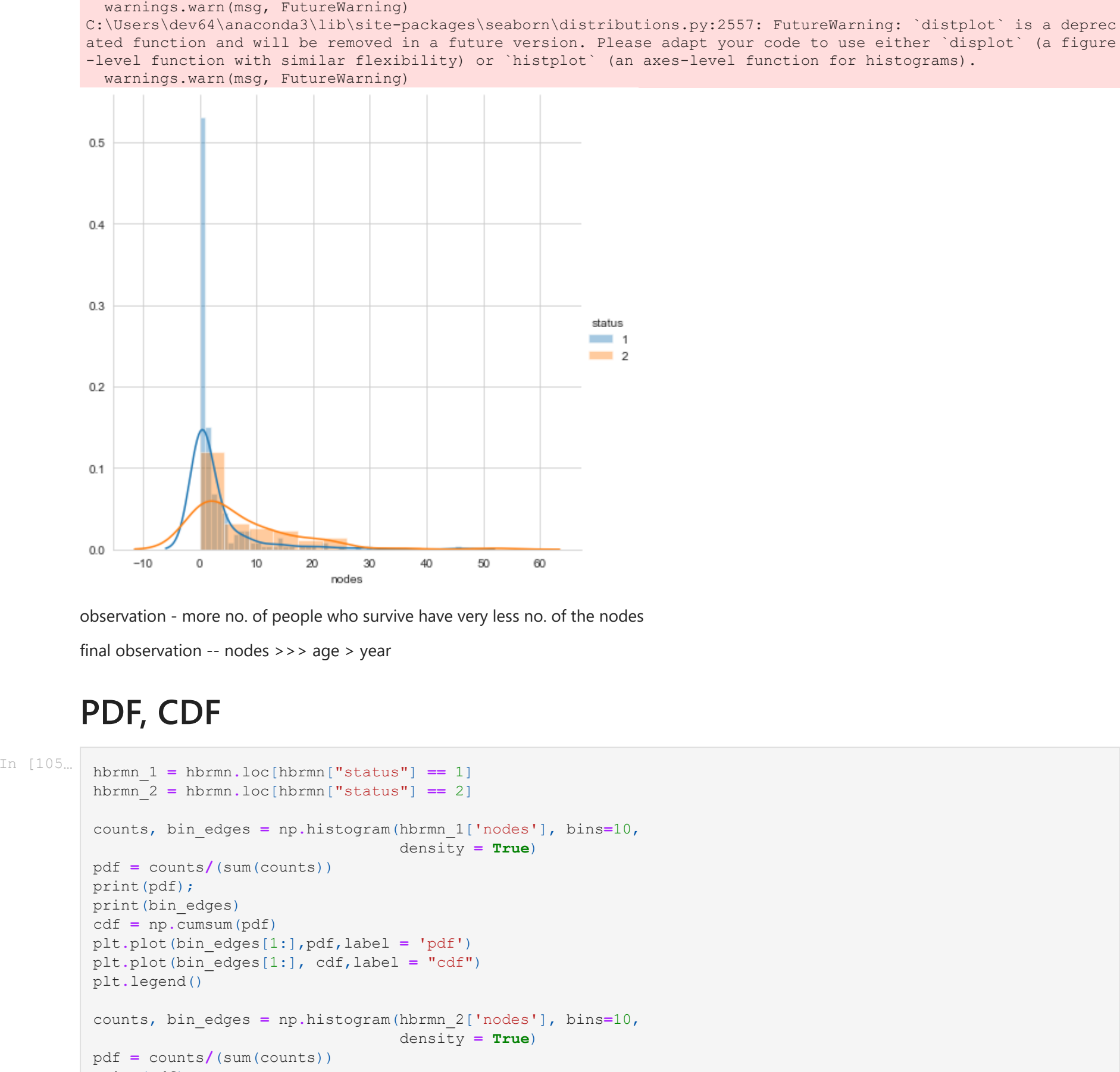
Histogram



observation - very hard to understand , because mostly data overlap.... even the pdf conside



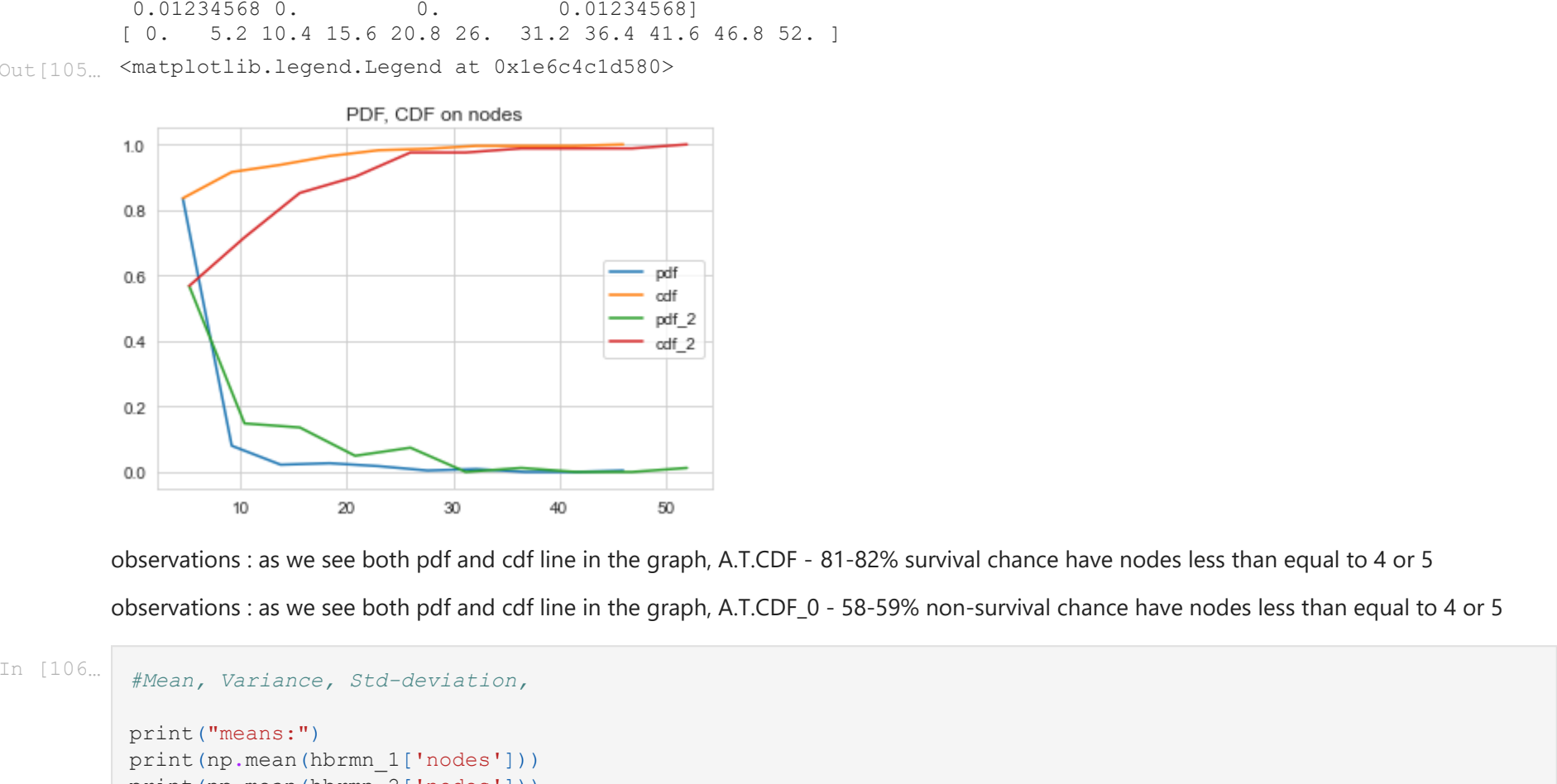
observation - very hard to understand , because mostly data overlap even pdf overlap



observation - more no. of people who survive have very less no. of the nodes

final observation -- nodes >>> age > year

PDF, CDF



observations : as we see both pdf and cdf line in the graph. A.T.CDF - 81-82% survival chance have nodes less than equal to 4 or 5
observations : as we see both pdf and cdf line in the graph. A.T.CDF,0 - 58-59% non-survival chance have nodes less than equal to 4 or 5

```
In [106]: #Mean, Variance, Std-deviation,

print ("means:")
print (np.mean(hbrmn_1['nodes']))
print (np.mean(hbrmn_2['nodes']))

print ("\nSTD-dev:");
print (np.std(hbrmn_1['nodes']))
print (np.std(hbrmn_2['nodes']))

Means:
2.7911111111111113
7.45679012345679

STD-dev:
5.857258449412131
9.1287760761632
```

```
In [107]: #Median, Quantiles, Percentiles, IQR.

print ("\nMedians:")
print (np.median(hbrmn_1['nodes']))
print (np.median(hbrmn_2['nodes']))

print ("\nQuantiles:")
print (np.percentile(hbrmn_1['nodes'],np.arange(0, 100, 25)))
print (np.percentile(hbrmn_2['nodes'],np.arange(0, 100, 25)))

print ("\n90th Percentiles:")
print (np.percentile(hbrmn_1['nodes'],90))
print (np.percentile(hbrmn_2['nodes'],90))

from statsmodels import robust
print ("robust Median Absolute Deviation")
print (robust.mad(hbrmn_1['nodes']))
print (robust.mad(hbrmn_2['nodes']))

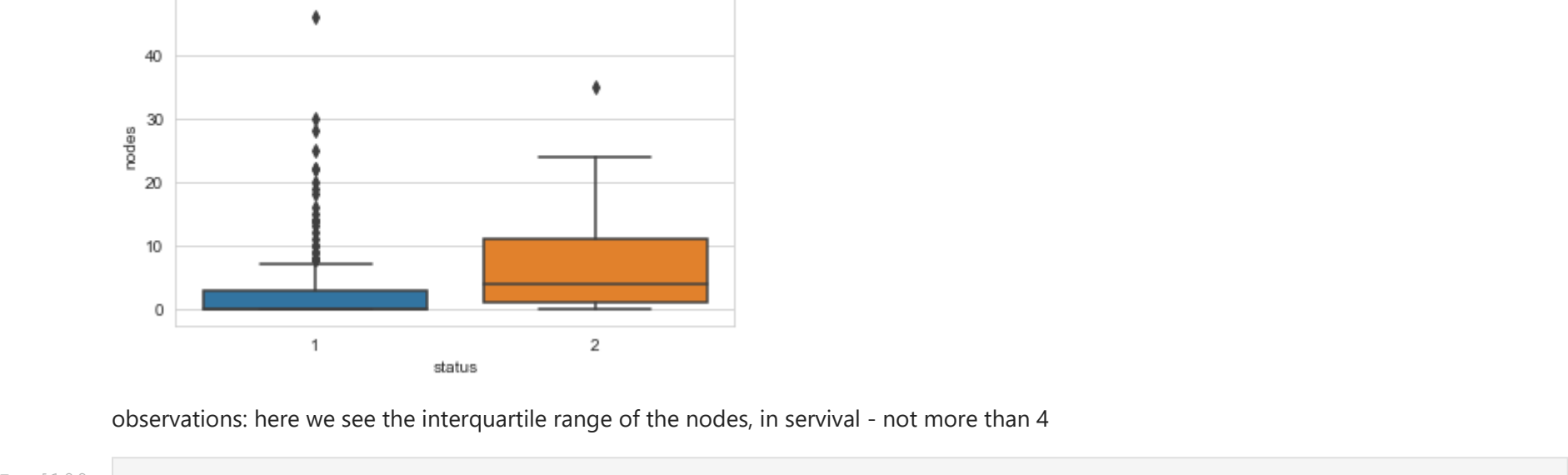
Medians:
0.0
4.0

Quantiles:
[0. 0. 0. 3.]
[0. 1. 4. 11.]

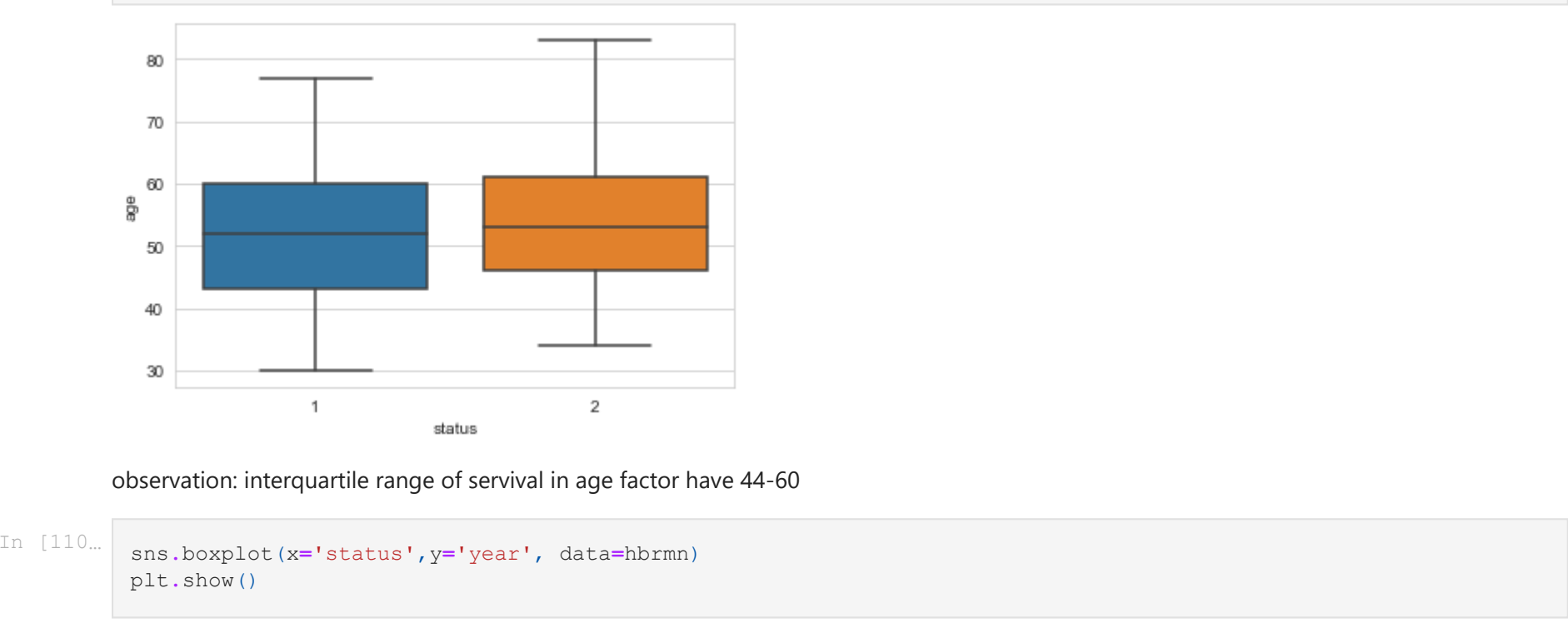
90th Percentiles:
8.0
20.0

Median Absolute Deviation
0.0
5.93040874022408
```

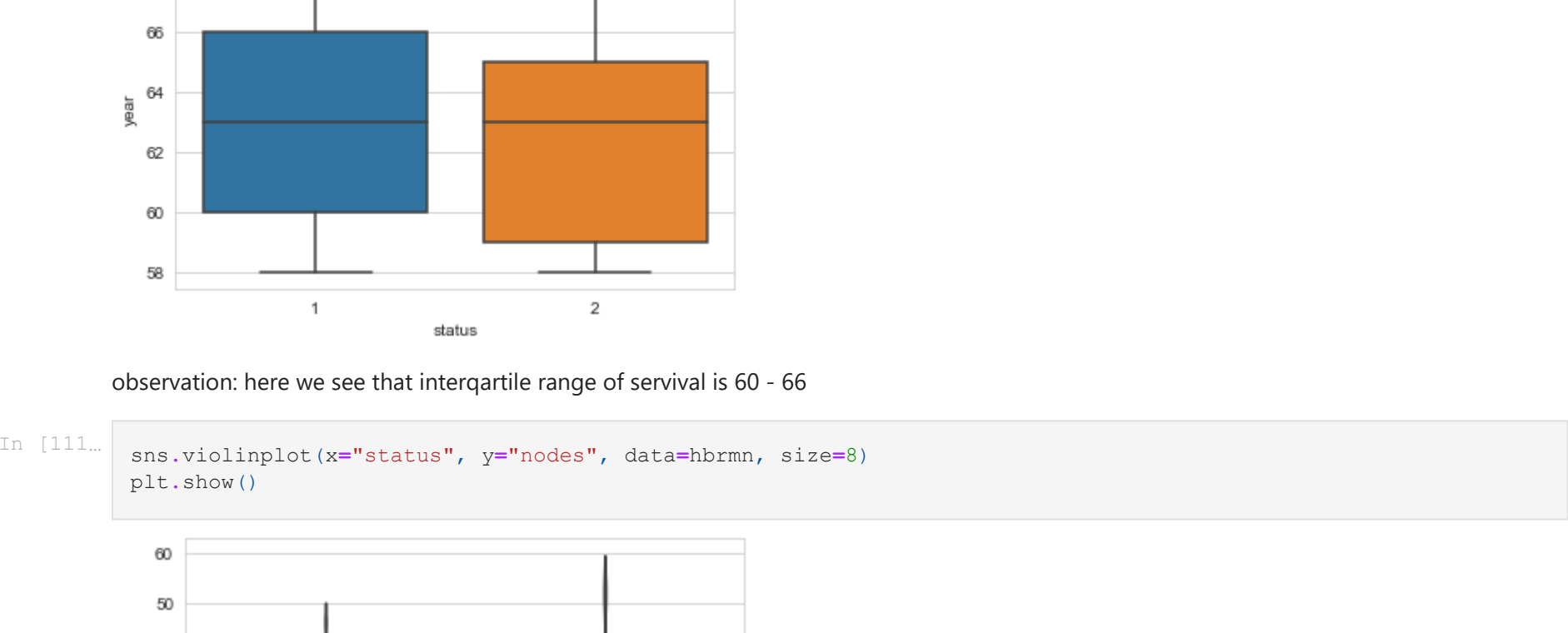
Box plot



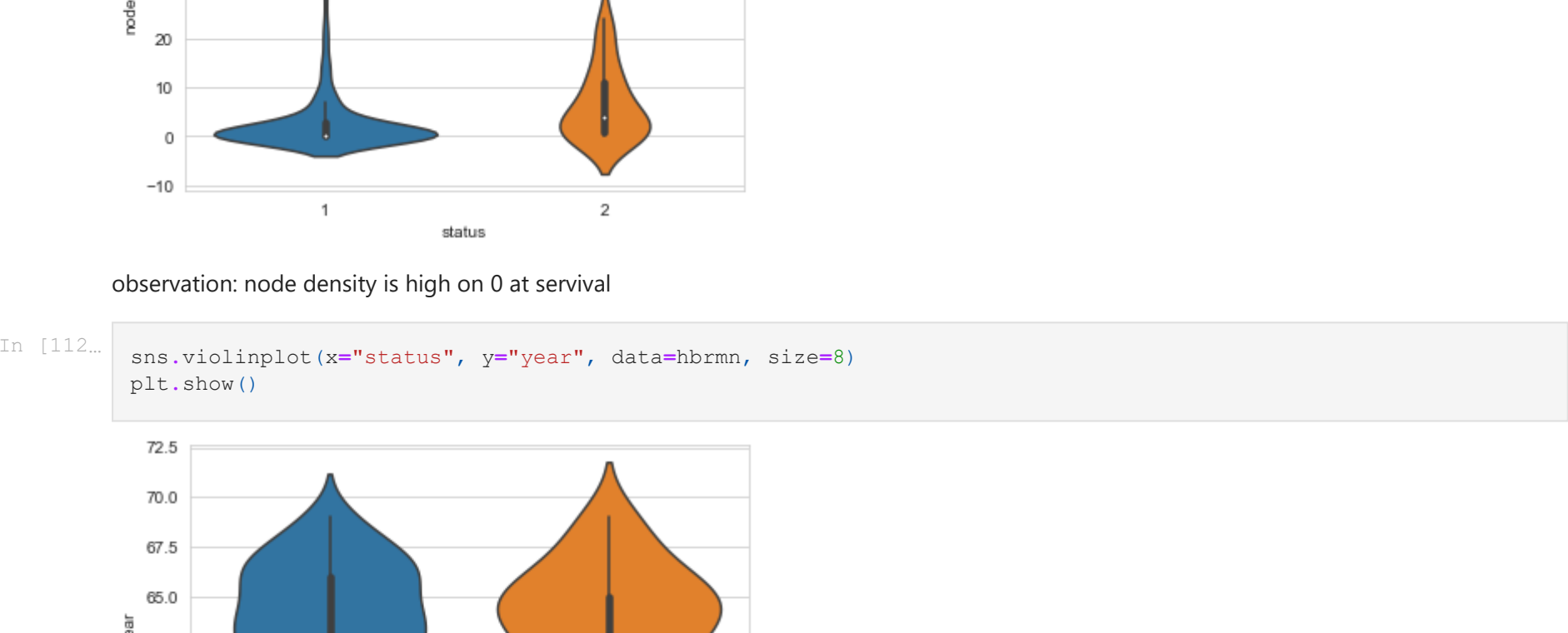
observations: here we see the interquartile range of the nodes, in survival - not more than 4



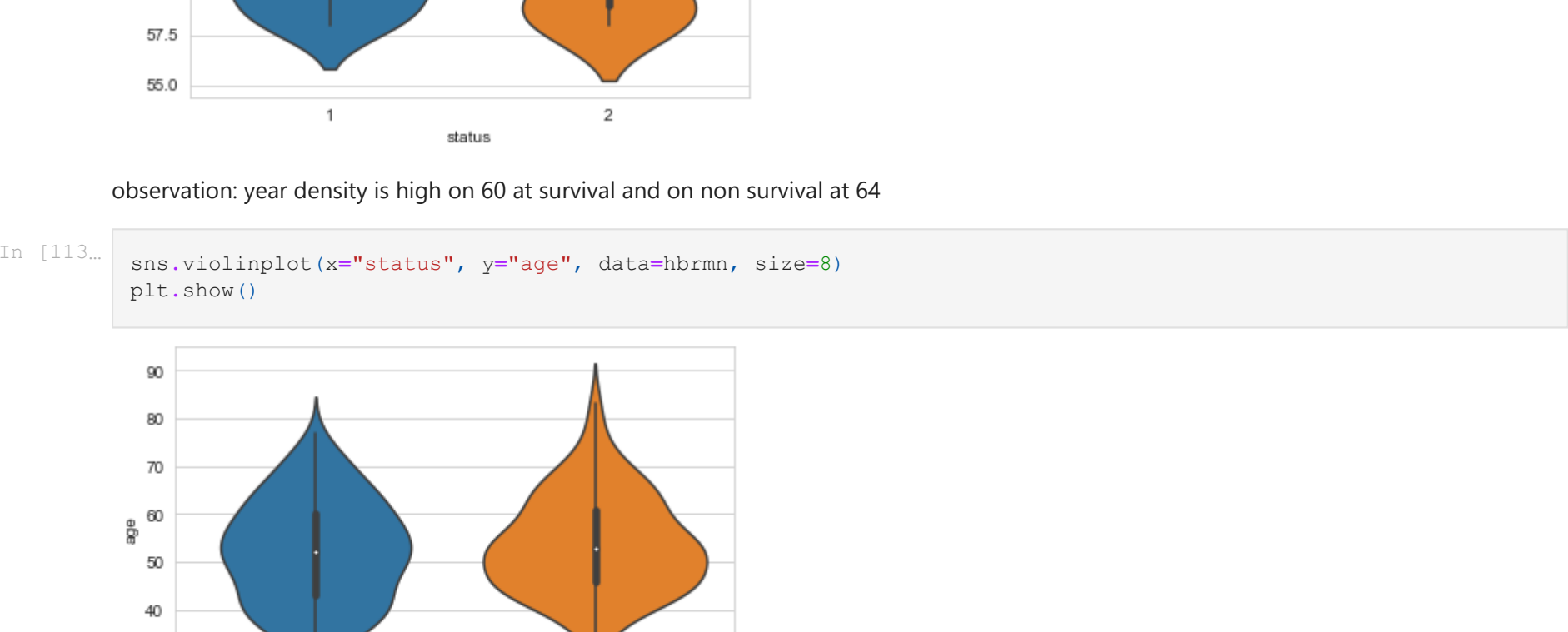
observation: interquartile range of survival in age factor have 44-60



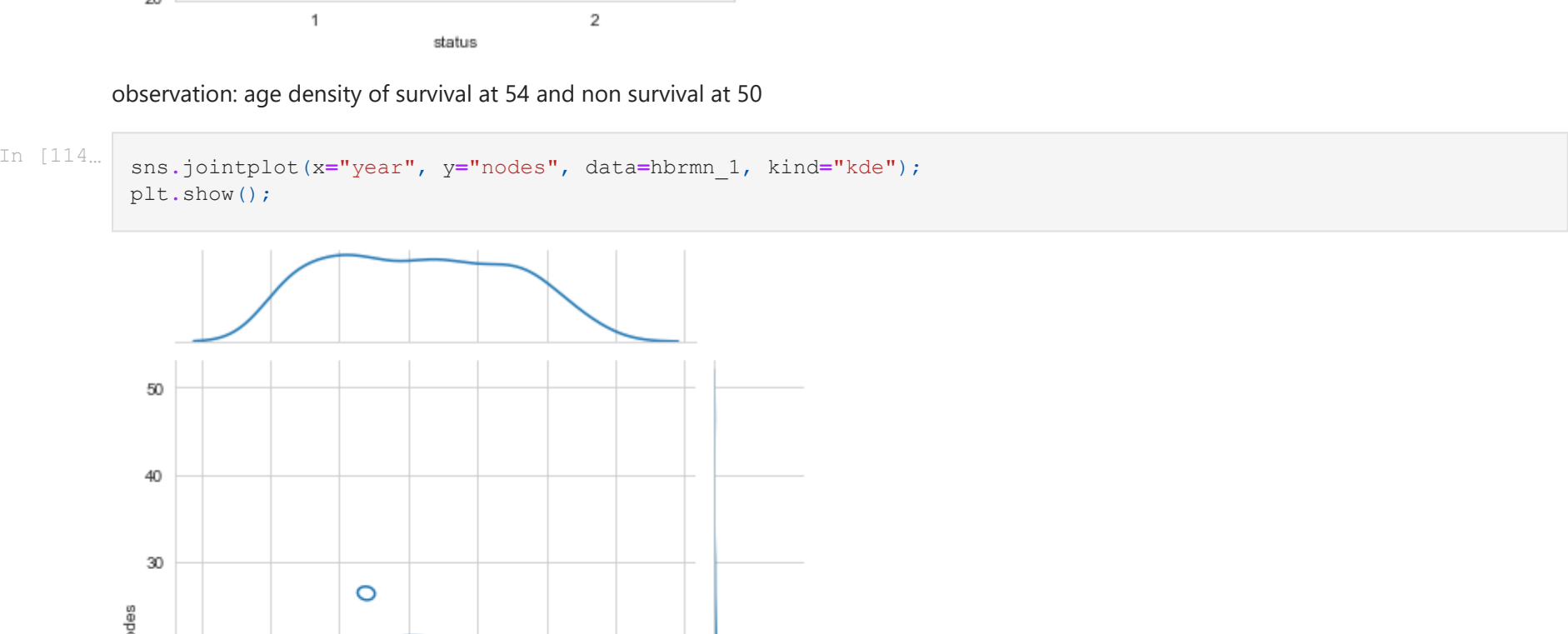
observation: here we see that interquartile range of survival is 60 - 66



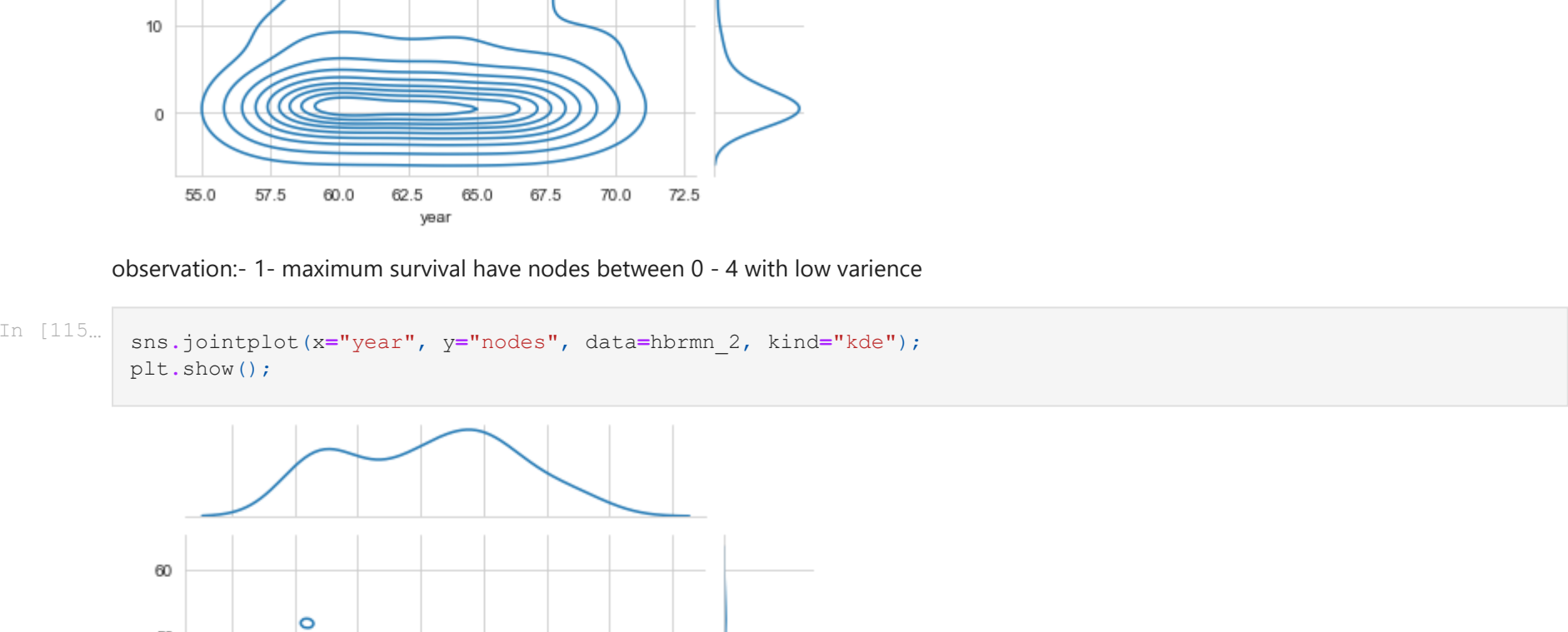
observation: node density is high on 0 at survival



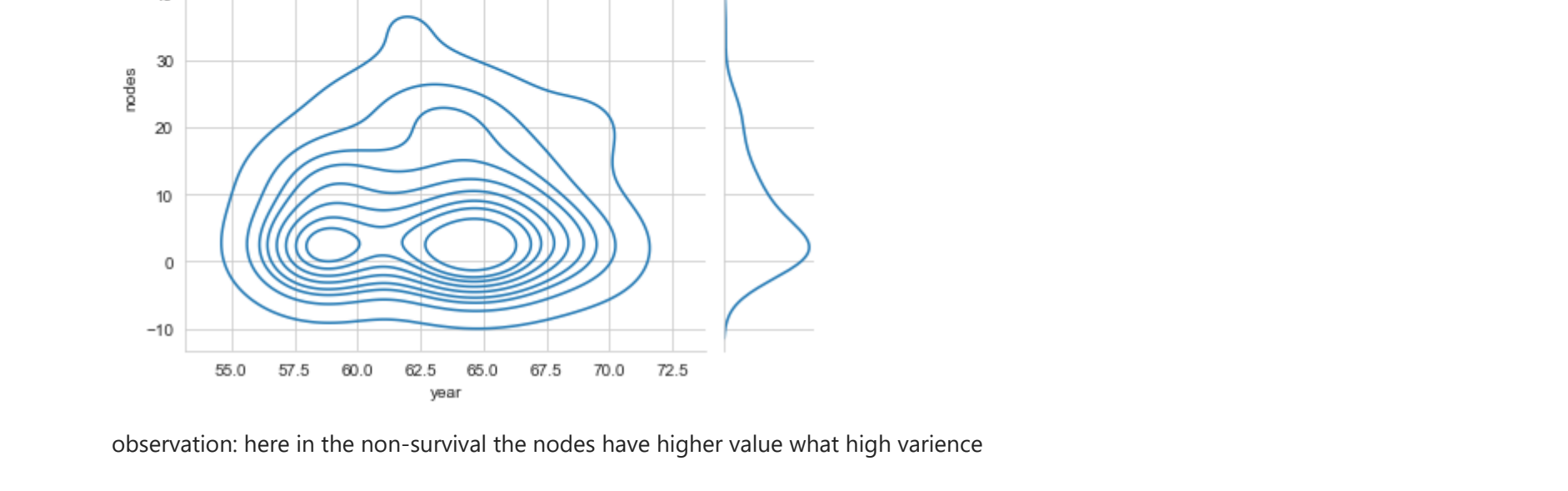
observation: year density is high on 60 at survival and on non survival at 64



observation: age density of survival at 54 and non survival at 50



observation:- 1- maximum survival have nodes between 0 - 4 with low variance



observation: here in the non-survival the nodes have higher value what high variance