

Exploratory data Analysis(EDA) for Haberman cancer survival status

Description:

The dataset contains cases from a study that was conducted between 1958 and 1970 at the University of Chicago's Billings Hospital on the survival of patients who had undergone surgery for breast cancer.

- Number of DataPoints: 306
- Number of Features: 4 (including the class)

Feature Information:

- Age of patient at time of operation (numerical)
- Patient's year of operation (year - 1900, numerical)
- Number of positive axillary nodes detected (numerical)
- Survival status (class)
 - 1(survived) = the patient survived 5 years or longer.
 - 2(died) = the patient died within 5 year.

1. High level statistics:

```
In [1]: import matplotlib.pyplot as plt
import numpy as np
import seaborn as sns
import pandas as pd

In [21]: haber = pd.read_csv('haberman.csv')

In [22]: print('Total Datapoints:',haber.shape[0],\
'\nTotal Features:',haber.shape[1],'\n')

print('Features:')
for val, col in enumerate(haber.columns):
    print(val+1,':',col)

Total datapoints: 306
Total Features: 4

Features:
1 : age
2 : year
3 : nodes
4 : status
```

Data-points per class:

```
In [23]: print('Number of classes:', haber['status'].unique().size,'\n')

#Mapping target attribute numeric values to string values
haber['status'] = haber['status'].map({ 1:'survived', 2:'died' })

print('DataPoints per class:')
print(haber['status'].value_counts())

Number of classes: 2

DataPoints per class:
survived    225
died         81
Name: status, dtype: int64
```

Observations:

1. Based on datapoints per class it's an imbalanced dataset.

2. Objective:

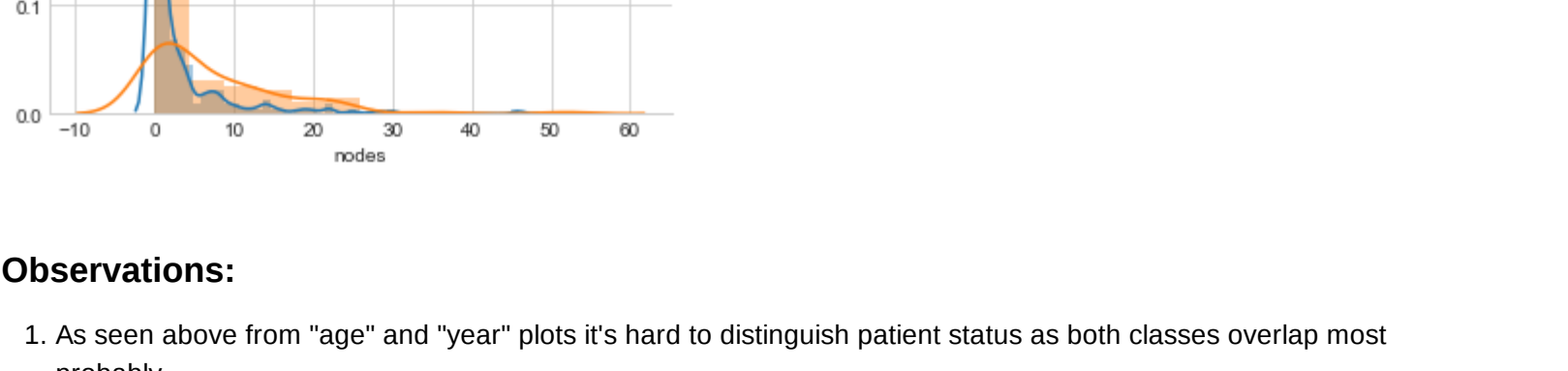
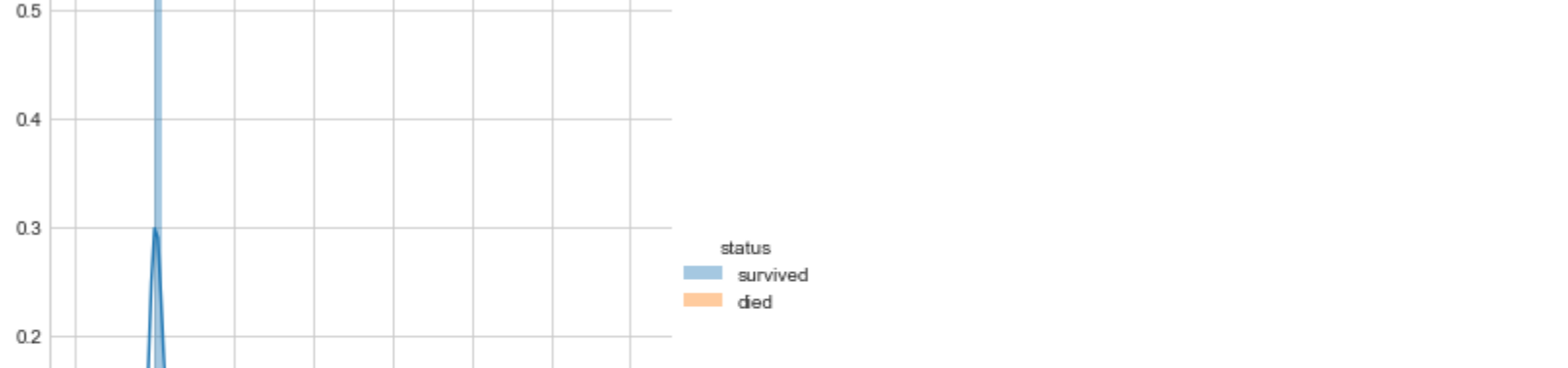
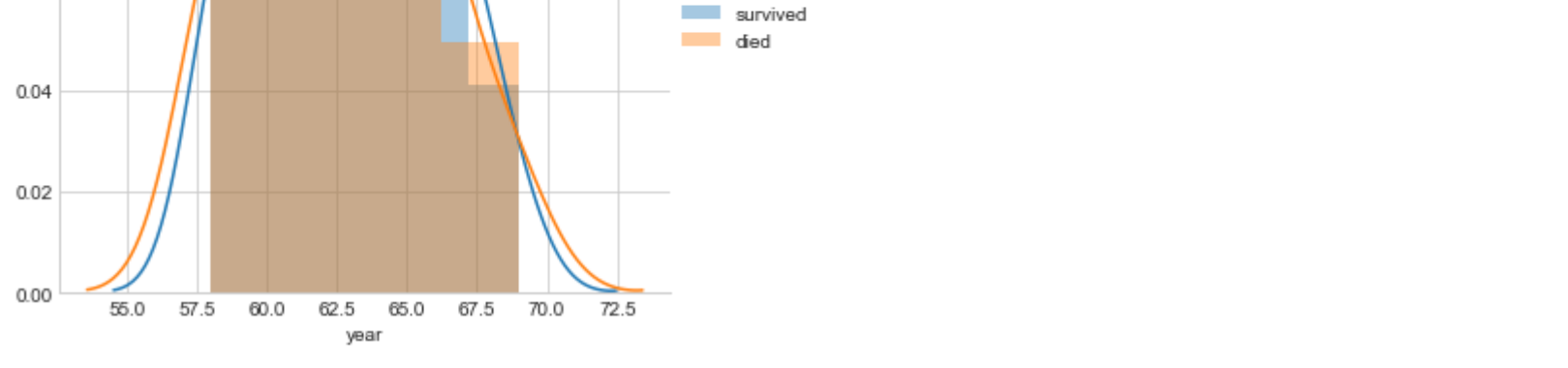
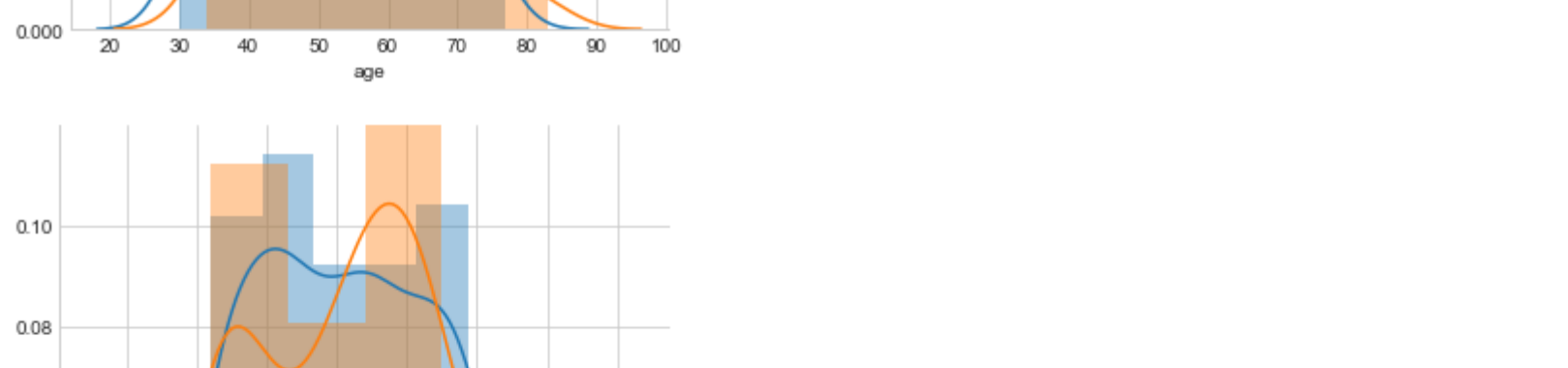
To classify whether a patient will survive 5 and more years or died with in 5 years from the year of treatment based on the given features:

- Age.
- Year of operation and
- Axillary lymph nodes.

3. Univariate analysis:

```
In [24]: import warnings
warnings.filterwarnings('ignore')

#PDF
sns.set_style('whitegrid')
for col in haber.columns[1:]:
    sns.FacetGrid(haber, hue='status', size=5)\
        .map(sns.distplot, col).add_legend()
plt.show()
plt.close()
```



Observations:

1. As seen above from "age" and "year" plots it's hard to distinguish patient status as both classes overlap most probably.
2. It can be seen that patients survived more than 5 years is dense from 0 to 5(approx) axillary "nodes". This feature which seems to be useful.

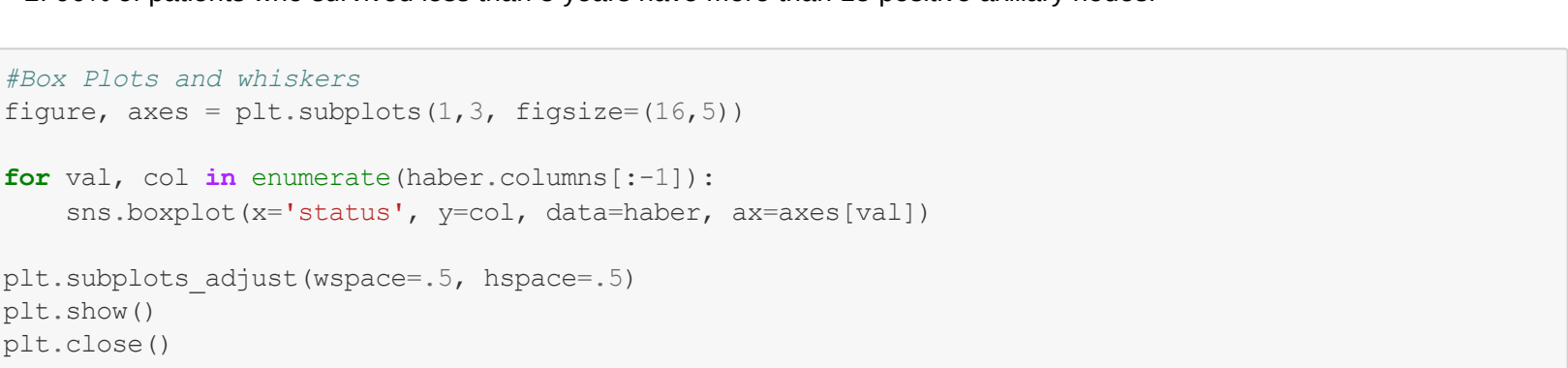
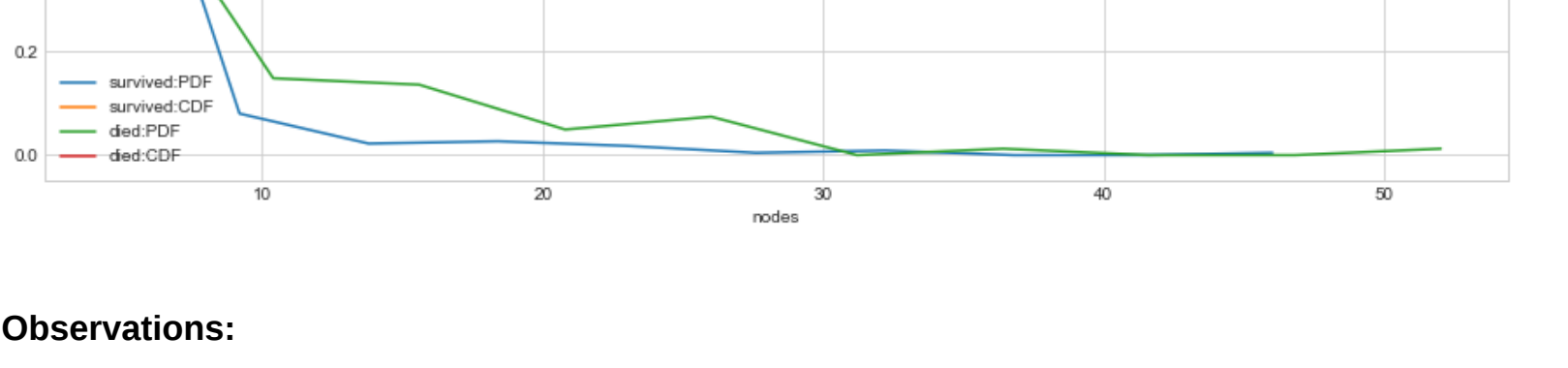
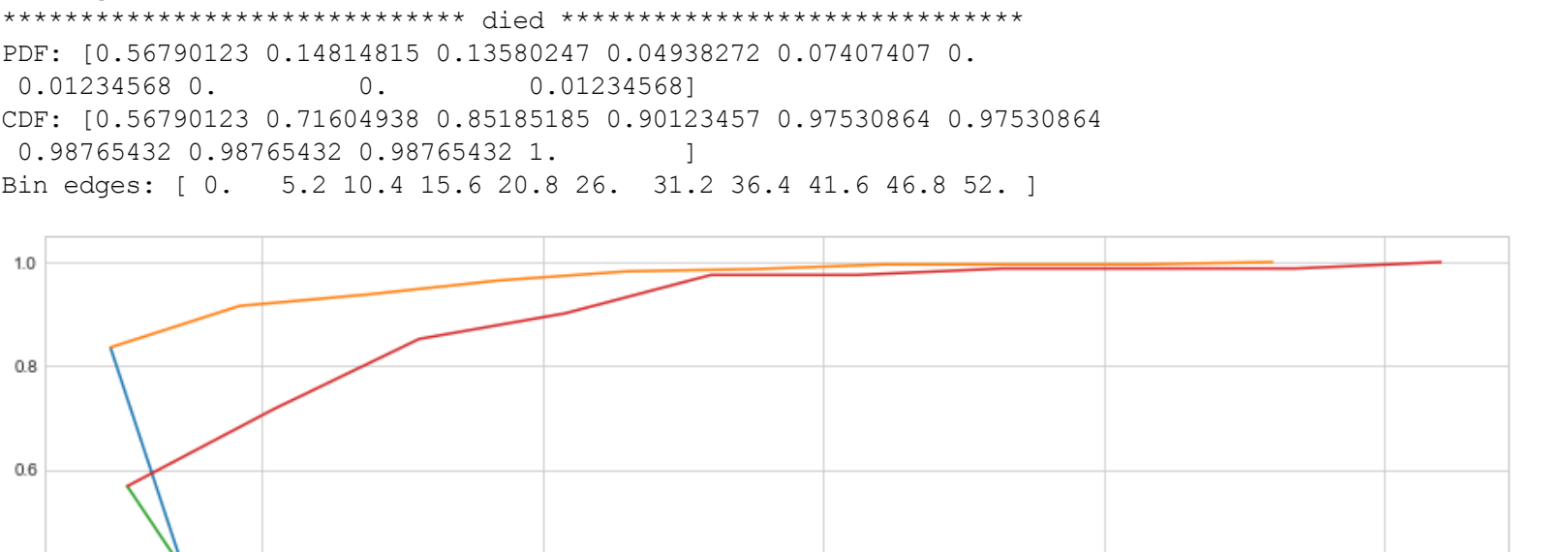
```
In [25]: #CDF
plt.figure(1, figsize=(15, 6))
for val in haber['status'].unique():
    counts, bin_edges = np.histogram(haber[ haber['status']==val ][ 'nodes' ],\
        counts, bin_edges = np.histogram bins=10, density = True)

    print('***30,val,***30)
    pdf = counts/sum(counts)
    cdf = np.cumsum(pdf)

    print('PDF:',pdf)
    print('CDF:',cdf)
    print('Bin edges:',bin_edges)

    plt.plot(bin_edges[1:], pdf, label='():PDF'.format(val))
    plt.plot(bin_edges[1:], cdf, label='():CDF'.format(val))

plt.legend()
plt.xlabel('nodes')
plt.show()
plt.close()
```



Observations:

1. 91% of patients who survived more than 5 years have less than 10 positive axillary nodes.
2. 90% of patients who survived less than 5 years have more than 15 positive axillary nodes.

```
In [26]: #Box Plots and whiskers
figure, axes = plt.subplots(1,3, figsize=(16,5))

for val, col in enumerate(haber.columns[1:]):
    sns.boxplot(x='status', y=col, data=haber, ax=axes[val])

plt.subplots_adjust(wspace=.5, hspace=.5)
plt.show()
plt.close()
```

