

## Exploratory Data Analysis

Exploratory Data Analysis One of the Important Part of Machine Learning, to understand How the feature are distributed, How the feature impact the model, To understand about outlier, mean, median, Iqr, Understand the data Based Visualization

Haberman's 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.

In this DataSet, Analysis How Many Patients are Survive More than Five Year after surgery, How many Patients are survive less than five years, Lets Start with our Analysis

```
In [2]: #import the necessary Library
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

```
In [3]: # read the Csv File store into the variable
df = pd.read_csv("%Content/Haberman.csv")
print(df.head())
```

age	year	nodes	status
0	30	64	1
1	30	62	3
2	30	65	0
3	31	59	2
4	31	65	4

## Observation

1, From this Dataset totally Four Features are present(Age, year, nodes and status)

2, Age, Year, nodes are Independent feature, Feature Status is Dependent feature (or) target feature (or) output feature

```
In [7]: print('Shape of the Given Data Set : ', df.shape)
Shape of the Given Data Set : (306, 4)
```

Totally 306 patients details are present in the Dataset, Each Patient have a details like Age, year of operation, No of lymph nodes and Survival Status Size of the Data - Nothing but how many data points are there in the data set. There are Four features, Each Feature has 306 data points, Total Size of our Given Data as 306 \* 4 = 1224 data points are used

```
In [9]: # print the what are the columns in the data and understand the type of each feature
for column in df.columns:
    print("columnName : {} and Type of the feature {}".format(column, type(df[column].dtypes)))
```

columnName : age and Type of the feature <class 'numpy.dtype'>  
columnName : year and Type of the feature <class 'numpy.dtype'>  
columnName : nodes and Type of the feature <class 'numpy.dtype'>  
columnName : status and Type of the feature <class 'numpy.dtype'>

```
In [12]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 306 entries, 0 to 305
Data columns (total 4 columns):
 #   Column  Non-Null Count  Dtype
---  --
 0   age     306 non-null       int64
 1   year    306 non-null       int64
 2   nodes   306 non-null       int64
 3   status  306 non-null       int64
dtypes: int64(4)
memory usage: 9.7 KB
```

We can see the snippet of the code explains given information of all features, Four Features in the dataset (Each feature have 306 data point) and there is Missing value and all the feature have a dtype integer

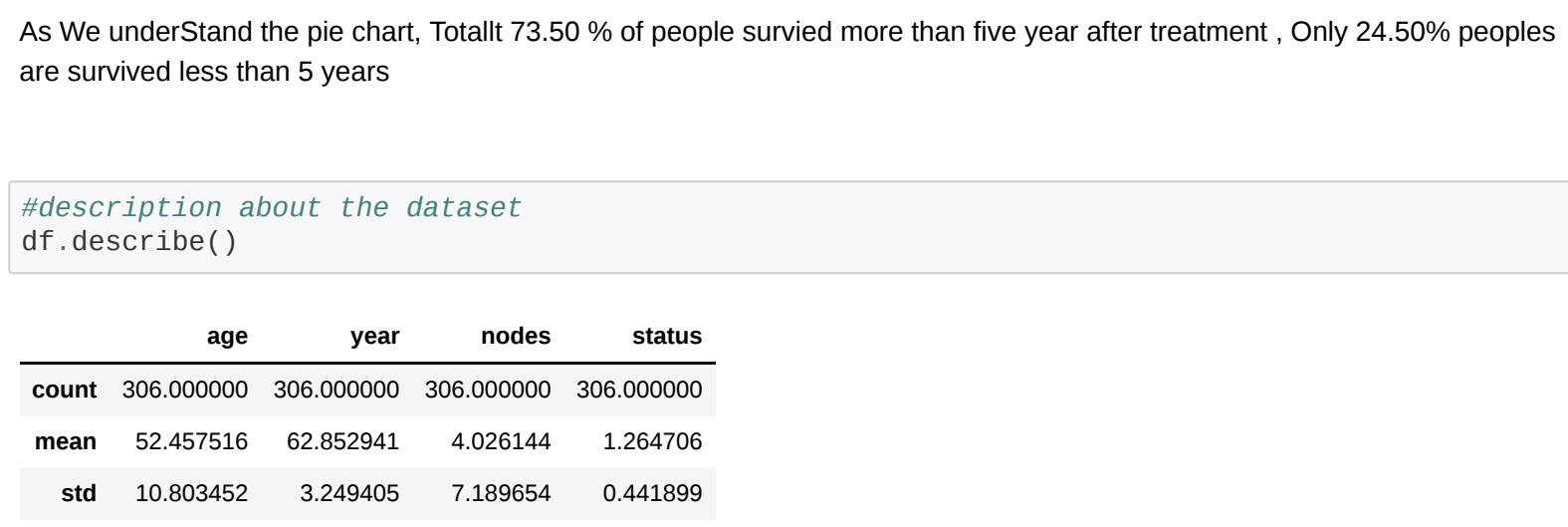
```
In [14]: df.groupby(['status']).count().T
```

```
Out[14]:
status    1    2
age      225  81
year     225  81
nodes    225  81
```

```
In [15]: df.status.value_counts()
```

```
Out[15]:
1    225
2     81
Name: status, dtype: int64
```

```
In [35]: plt.figure(figsize=(10,10))
plt.title("Survival Status of the patient")
labels = ["Survival Status", "Survival_Status(No)"]
colors = ["lightskyblue", "gold"]
df.status.value_counts().plot(kind='pie', explode = [0.1,0], autopct = '%0.2f%%', shadow = True, label
s = labels, colors = colors)
plt.legend(labels, loc = 'upper right')
plt.show()
```



As We understand the pie chart, Total 73.50 % of people survived more than five year after treatment , Only 24.50% peoples are survived less than 5 years

```
In [184]: #description about the dataset
df.describe()
```

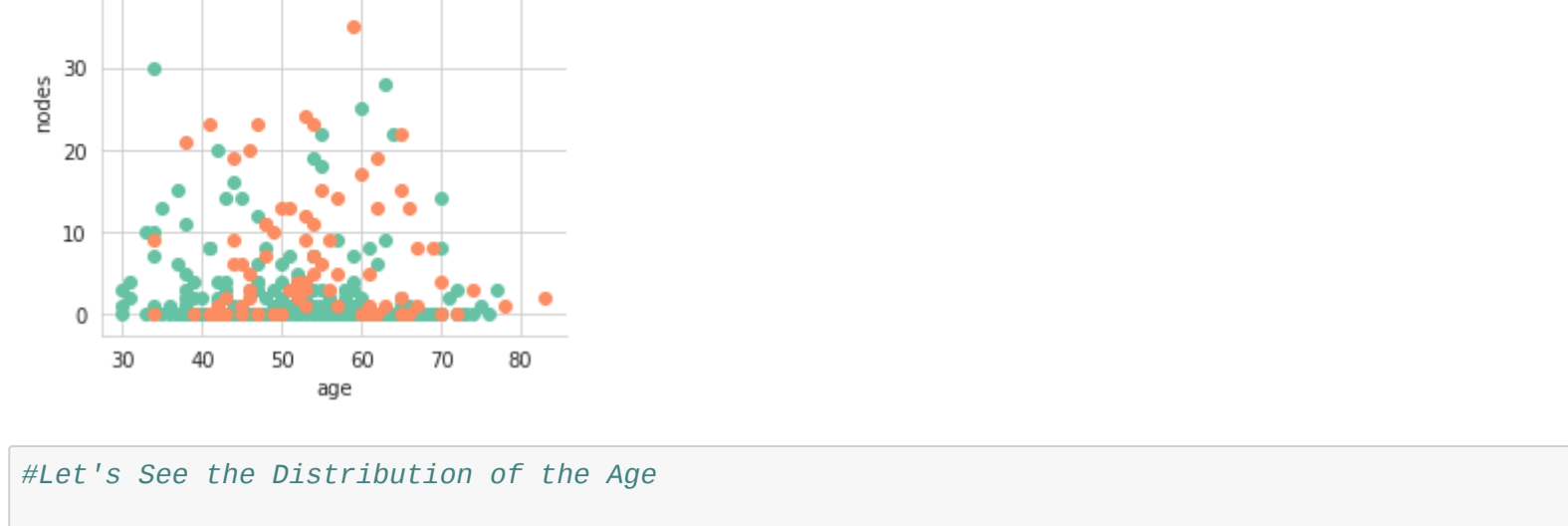
```
Out[184]:
```

	age	year	nodes	status
count	306.000000	306.000000	306.000000	306.000000
mean	52.47969	62.86204	4.00614	1.346706
std	10.892462	3.248405	7.189654	0.431099
min	30.000000	58.000000	0.000000	1.000000
25%	44.000000	60.000000	0.000000	1.000000
50%	52.000000	63.000000	1.000000	1.000000
75%	60.750000	65.750000	4.000000	2.000000
max	83.000000	69.000000	52.000000	2.000000

This shows the Age of the people less than 30 to 83 Although the maximum number of positive axillary nodes observed is 52, nearly 75% of the patients have less than 5 positive axillary nodes and nearly 25% of the patients have no positive axillary nodes

```
In [41]: plt.figure(figsize=(30,10))
sns.countplot(x = 'nodes', hue = 'status', data = df, palette='Set3')
plt.title('Survival Status for Patient on Feature Nodes')
```

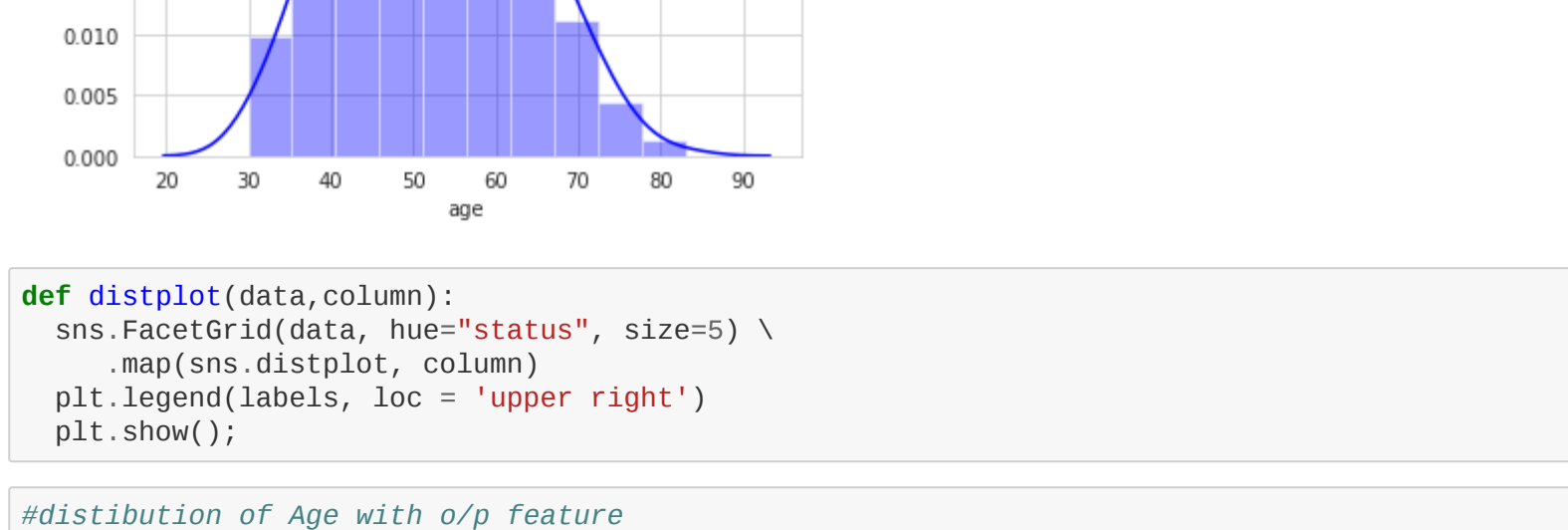
```
Out[41]: Text(0.5, 1.0, 'Survival Status for Patient on Feature Nodes')
```



This Figure shows that Nearly 120 no of patient Survived More than Five Year has no positive axillary nodes

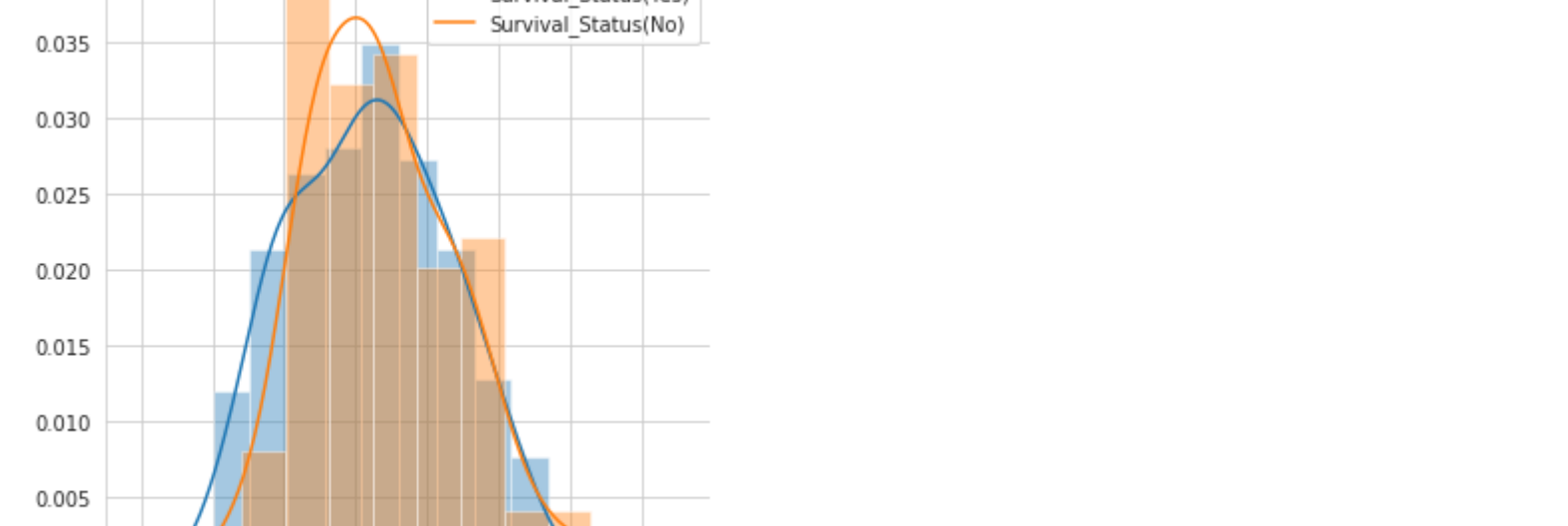
```
In [48]: sns.set_style('whitegrid')
plt.figure(figsize=(10,15))
sns.factorplot(df, hue='status', size=4, palette='Set2').map(plt.scatter, "age", "nodes")
plt.legend(labels, loc = 'upper right')
plt.show()
```

/usr/local/lib/python3.7/dist-packages/seaborn/axisgrid.py:316: UserWarning: The 'size' parameter has been renamed to 'height'; please update your code.
warnings.warn(msg, UserWarning)



```
In [53]: #Let's See the Distribution of the Age
sns.distplot(df['age'], bins = 10, color = 'blue')
plt.title('Distribution of Age')
```

```
Out[53]: Text(0.5, 1.0, 'Distribution of Age')
```



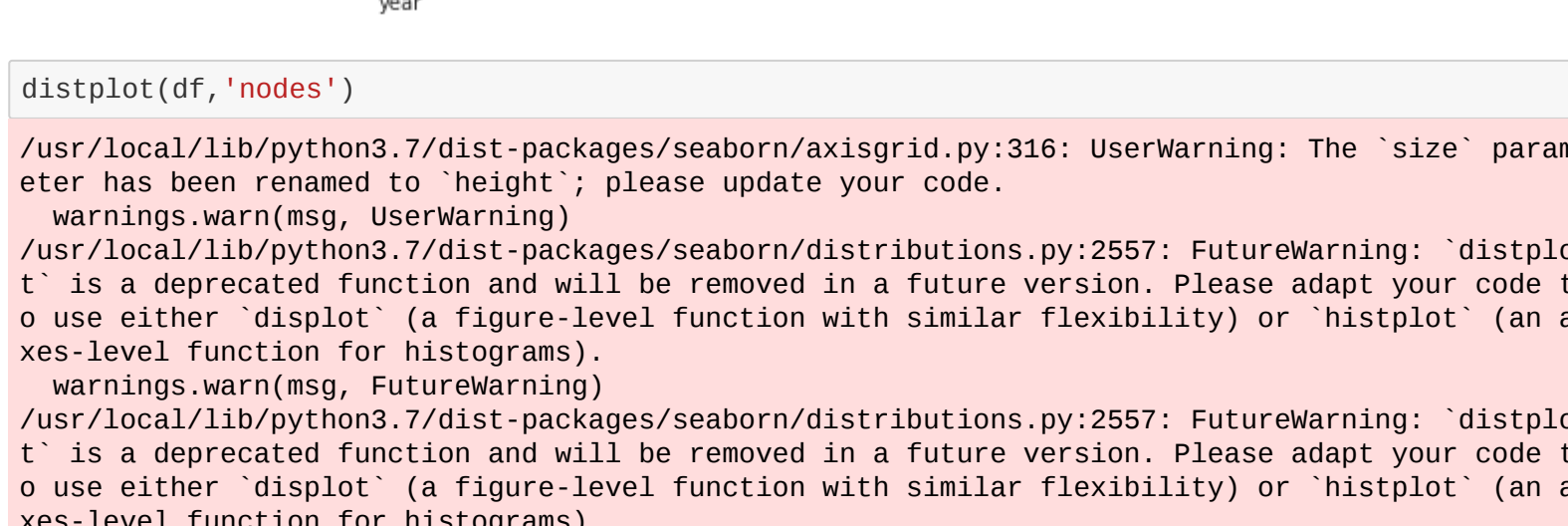
```
In [59]: def distplot(data, column):
sns.FacetGrid(data, hue='status', size=5) \
    .map(sns.distplot, column)
plt.legend(labels, loc = 'upper right')
plt.show()
```

```
In [60]: #distribution of year with or/p feature
distplot(df, 'age')
```

/usr/local/lib/python3.7/dist-packages/seaborn/axisgrid.py:316: UserWarning: The 'size' parameter has been renamed to 'height'; please update your code.
warnings.warn(msg, UserWarning)

/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2557: 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).

/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2557: 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).

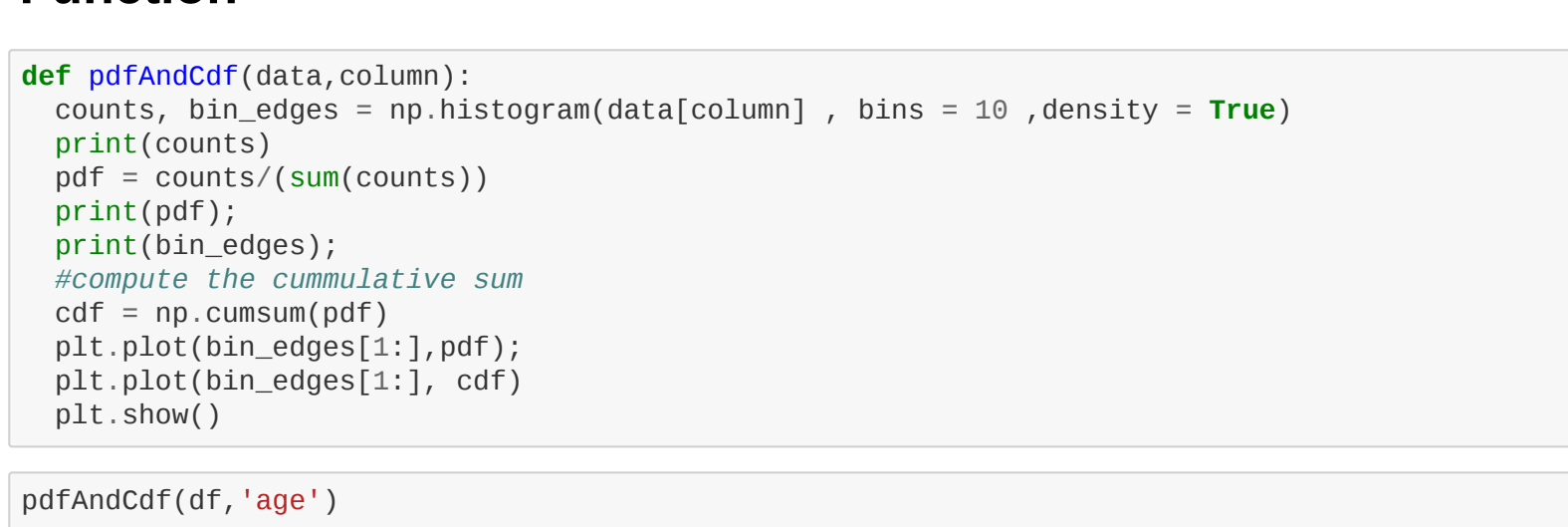


```
In [62]: #distribution of year with survival status
distplot(df, 'year')
```

/usr/local/lib/python3.7/dist-packages/seaborn/axisgrid.py:316: UserWarning: The 'size' parameter has been renamed to 'height'; please update your code.
warnings.warn(msg, UserWarning)

/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2557: 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).

/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2557: 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).

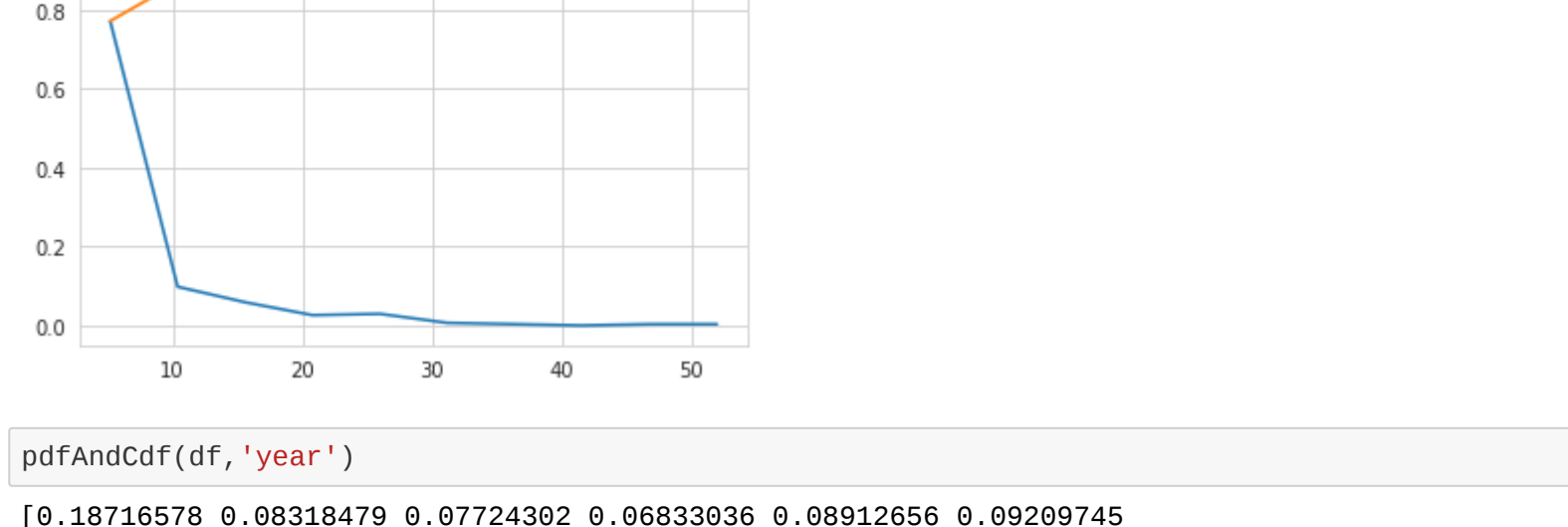


```
In [63]: distplot(df, 'nodes')
```

/usr/local/lib/python3.7/dist-packages/seaborn/axisgrid.py:316: UserWarning: The 'size' parameter has been renamed to 'height'; please update your code.
warnings.warn(msg, UserWarning)

/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2557: 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).

/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2557: 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).



From the univariate distribution, we can see that Positive\_Axillary\_Nodes is important feature that determines our dependent variable Status. The second important feature is Age. Observation: The number of positive lymph nodes of the survivors is highly denser from 0 to 5.

## Probability Density Function and Cumulative Density Function

```
In [69]: def pdfAndcdf(data, column):
count, bin_edges = np.histogram(data[column], bins = 10, density = True)
print(counts)
pdf = counts/sum(counts)
print(pdf)
#compute the cumulative sum
cdf = np.cumsum(pdf)
plt.plot(bin_edges[1:], pdf, color='green')
plt.plot(bin_edges[1:], cdf, color='red')
plt.show()
```

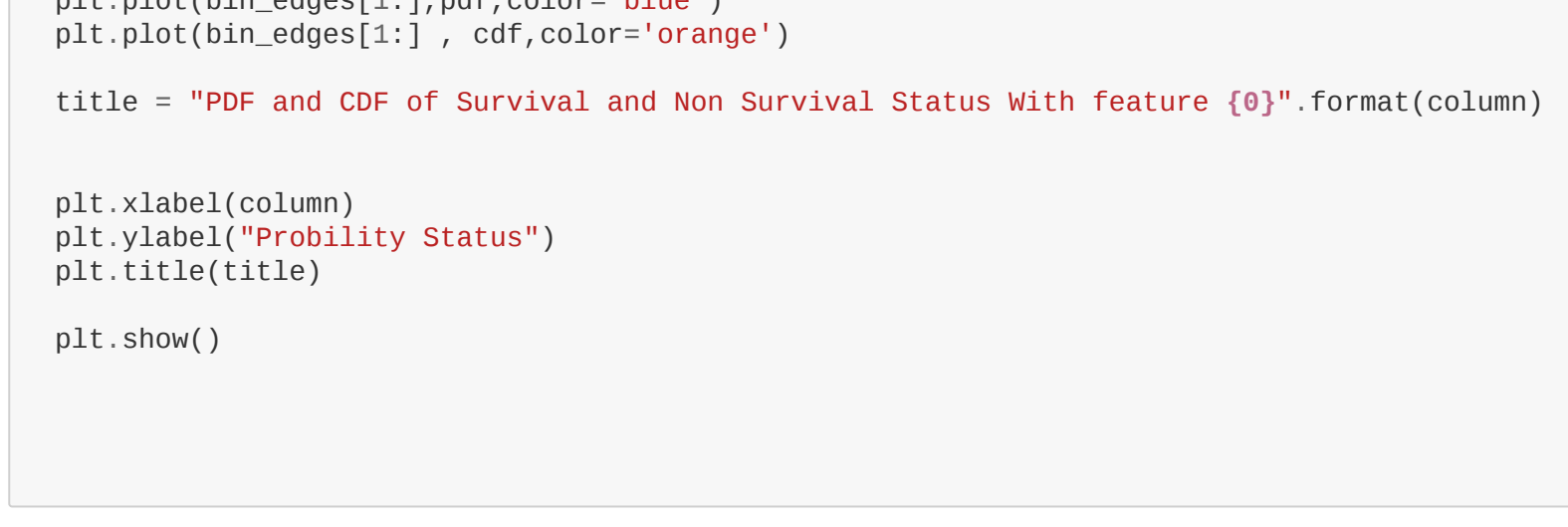
```
In [70]: pdfAndcdf(df, 'age')
```

```
Out[70]:
[0.00986558 0.01644817 0.02836355 0.03267974 0.03391294 0.02528855
0.02528855 0.01111111 0.07427213 0.07130125]
[0.02528855 0.08982322 0.09555556 0.17973856 0.09803822 0.10130719
0.09803822 0.09158327 0.06498732 0.07516354 0.09803822 0.07843137]
[0.09158327 0.09158327 0.08169935 0.08635955]
[0. 0. 5.2 18.4 45.9 51.2 56.5 61.8 67.1 72.4 77.7 83. ]
```



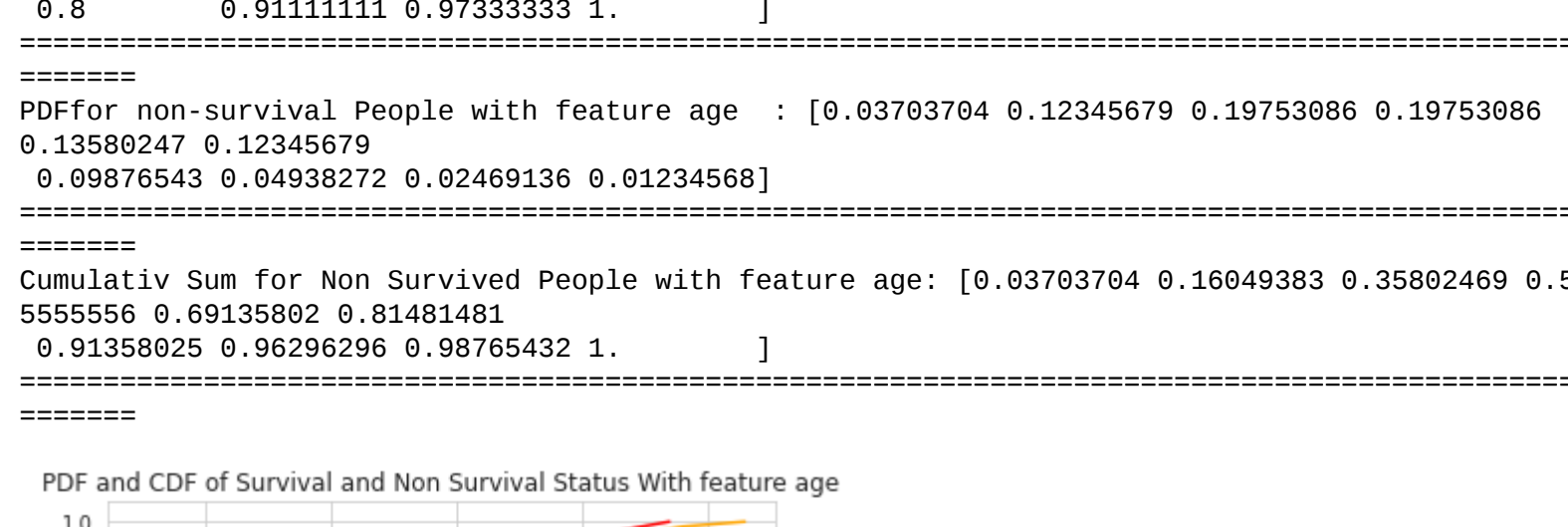
```
In [71]: pdfAndcdf(df, 'nodes')
```

```
Out[71]:
[0.14831574 0.0188537 0.01131222 0.0562765 0.09565561 0.09125691
0.09062846 0.09062846 0.09062846]
[0.09062846 0.09062846 0.09062846 0.09062846 0.09062846 0.09062846
0.09062846 0.09062846 0.09062846 0.09062846 0.09062846 0.09062846]
[0. 0. 5.2 18.4 45.9 51.2 56.5 61.8 67.1 72.4 77.7 83. ]
```



```
In [72]: pdfAndcdf(df, 'year')
```

```
Out[72]:
[0.18716578 0.08318479 0.0724302 0.06833036 0.08912556 0.09209745
0.08318479 0.08318479 0.0431619 0.07130125]
[0.08318479 0.08318479 0.08318479 0.08318479 0.08318479 0.08318479
0.08318479 0.08318479 0.08318479 0.08318479 0.08318479 0.08318479]
[0.08318479 0.08318479 0.08318479 0.08318479 0.08318479 0.08318479
0.08318479 0.08318479 0.08318479 0.08318479 0.08318479 0.08318479]
[0. 0. 5.2 18.4 45.9 51.2 56.5 61.8 67.1 72.4 77.7 83. ]
```



From the above plot Blue line Indicates PDF and Orange Line Indicates CDF we can see that almost 80% of the patients have positive lymph nodes less than 10.

```
In [92]: #plot Probability Density Function and CDF with Class Label
def PDFandCDFwithClassLabel(data, column):
#survived class label data
survived = data.loc[data['status'] == 1, :]
#non-survived
non_survived = data.loc[data['status'] == 2, :]
```

```
#plot for survived status
count, bin_edges = np.histogram(survived[column], bins = 10, density = True)
pdf = count/sum(count)
```

```
print("PDF of Survival Status with Feature {} is : {}".format(pdf, column))
```

```
#compute cdf
cdf = np.cumsum(pdf)
print("Cumulative Sum for Survived People with Feature {} is : {}".format(cdf, column))
```

```
plt.plot(bin_edges[1:], pdf, color='green')
plt.plot(bin_edges[1:], cdf, color='red')
print("==" * 100)
```

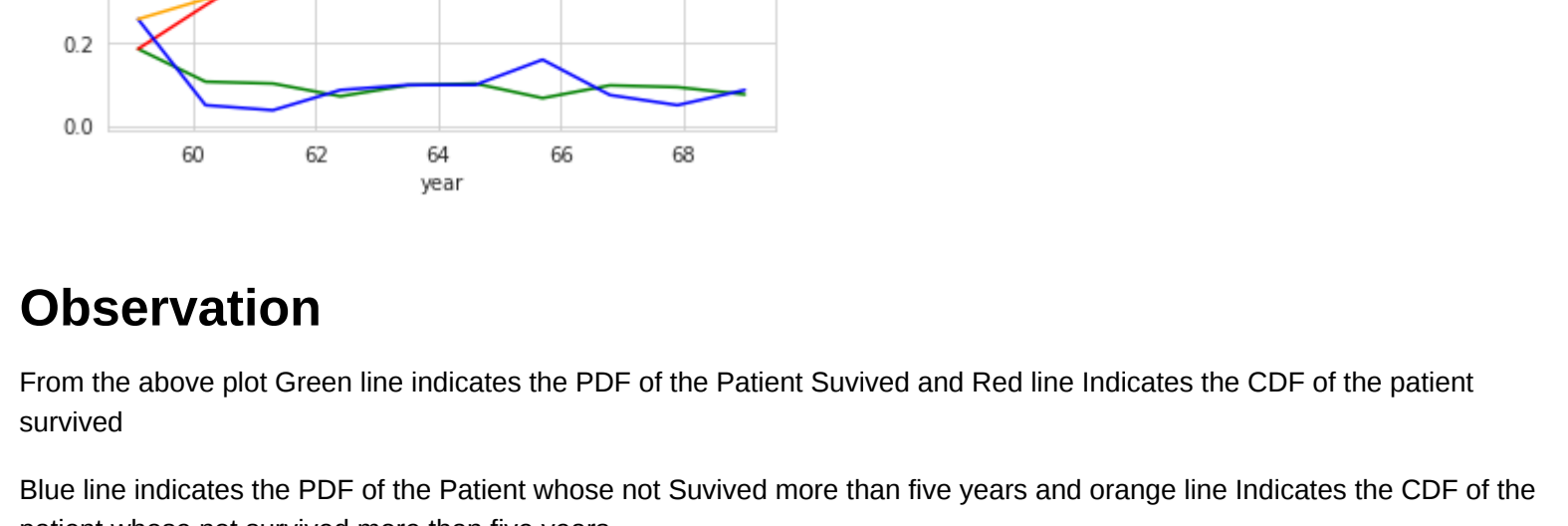
```
count, bin_edges = np.histogram(non_survived[column], bins = 10, density = True)
pdf = count/sum(count)
```

```
print("PDF for non-survived People with Feature {} is : {}".format(pdf, column))
```

```
#compute cdf
print("==" * 100)
cdf = np.cumsum(pdf)
print("Cumulative Sum for Non Survived People with Feature {} is : {}".format(cdf, column))
```

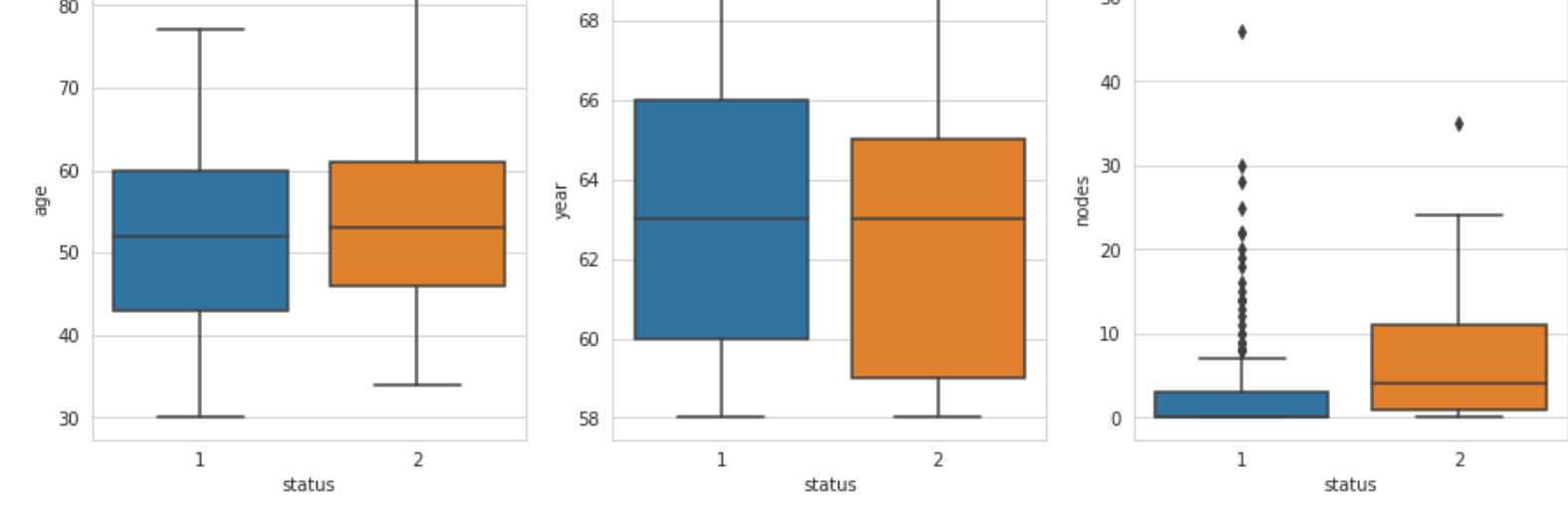
```
print("==" * 100)
plt.plot(bin_edges[1:], pdf, color='blue')
plt.plot(bin_edges[1:], cdf, color='orange')
```

```
title = "PDF and CDF of Survival and Non Survival Status With Feature {}".format(column)
plt.xlabel(column)
plt.ylabel('Probability Status')
plt.title(title)
plt.show()
```



```
In [94]: PDFandCDFwithClassLabel(df, 'nodes')
```

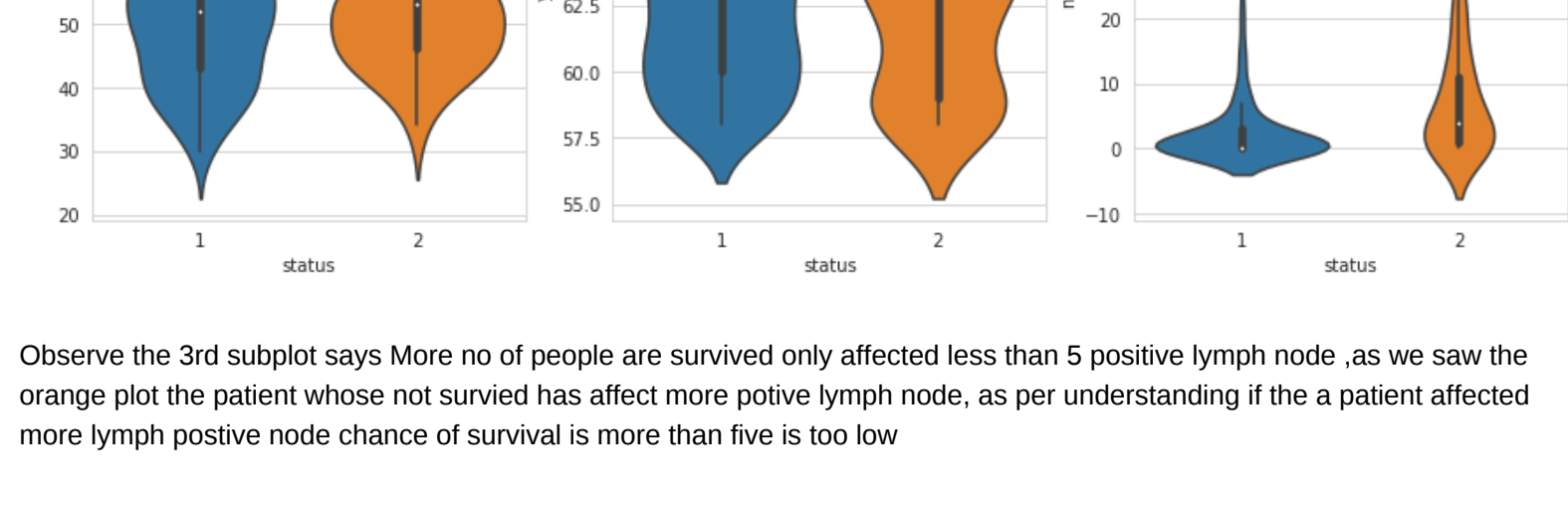
```
Out[94]:
PDF of Survival Status with feature age is : [0.05333333 0.16666667 0.12444444 0.09333333 0.16444444 0.16444444
0.09333333 0.11111111 0.06222222 0.02666667]
Cumulative Sum for Survived People with feature age : [0.05333333 0.16 0.28444444 0.37777778 0.54222222 0.70666667 0.83333333 0.94444444 1.0]
PDF for non-survived People with feature age : [0.0376704 0.12345679 0.19753086 0.3398247 0.32345679
0.09876543 0.04938272 0.02469136 0.01234568]
Cumulative Sum for Non Survived People with Feature age : [0.0376704 0.0376704 0.0376704 0.0376704 0.0376704 0.0376704
0.0376704 0.0376704 0.0376704 0.0376704 0.0376704 0.0376704]
PDF and CDF of Survival and Non Survival Status With feature age
```



Observe the 3rd subplot See More no of people are survived only affected less than 5 positive lymph node, As we saw the orange plot the patient whose not survived has affect more positive lymph node, as per understanding if the a patient affected more lymph positive node chance of survival is more than five is too low

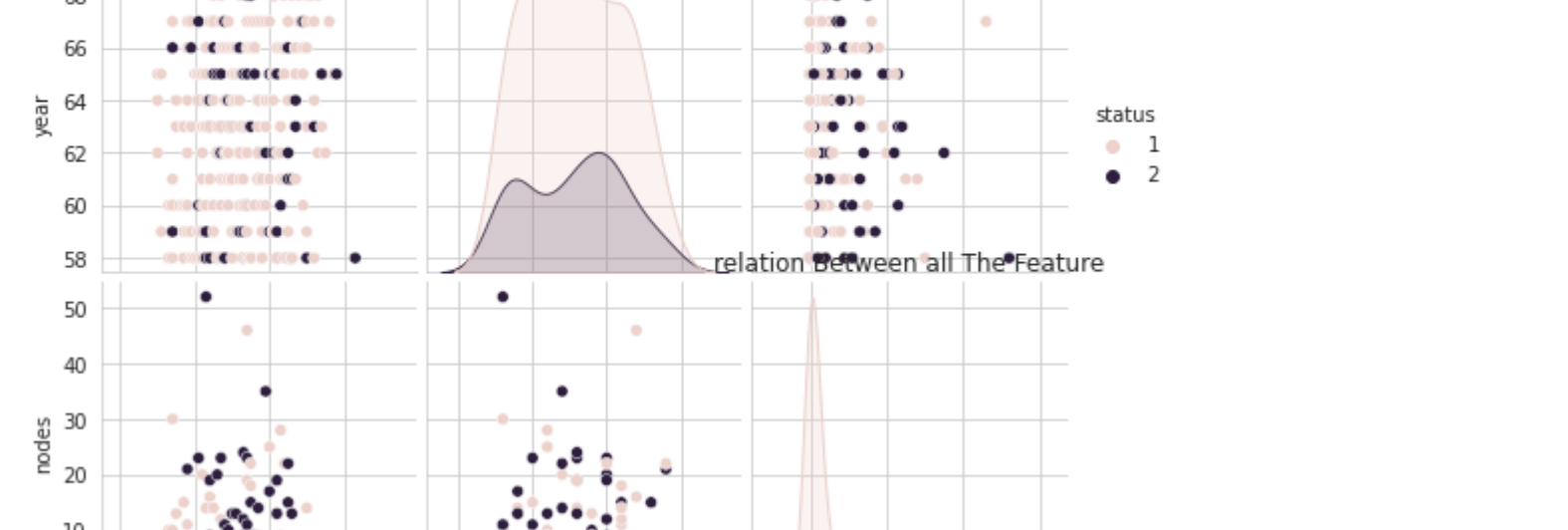
## Multivariate Analysis

```
In [97]: sns.pairplot(df, hue='status')
plt.title('relation Between all The Feature')
plt.show()
```



```
In [98]: #Bivariate Analysis
#Relationship Between Age and Year
sns.bivariateAnalysisScatterPlot(data, column1, column2):
#firstparam = columnname
#secondparam = columnname
#thirdparam = columnname
plt.figure(figsize=(15,10))
sns.set_style('whitegrid')
title = "Relationship Between {} and {}".format(column1, column2)
x_axis = column1
y_axis = column2
plt.title(title)
plt.xlabel(x_axis)
plt.ylabel(y_axis)
plt.scatter(df[column1], df[column2])
plt.show()
```

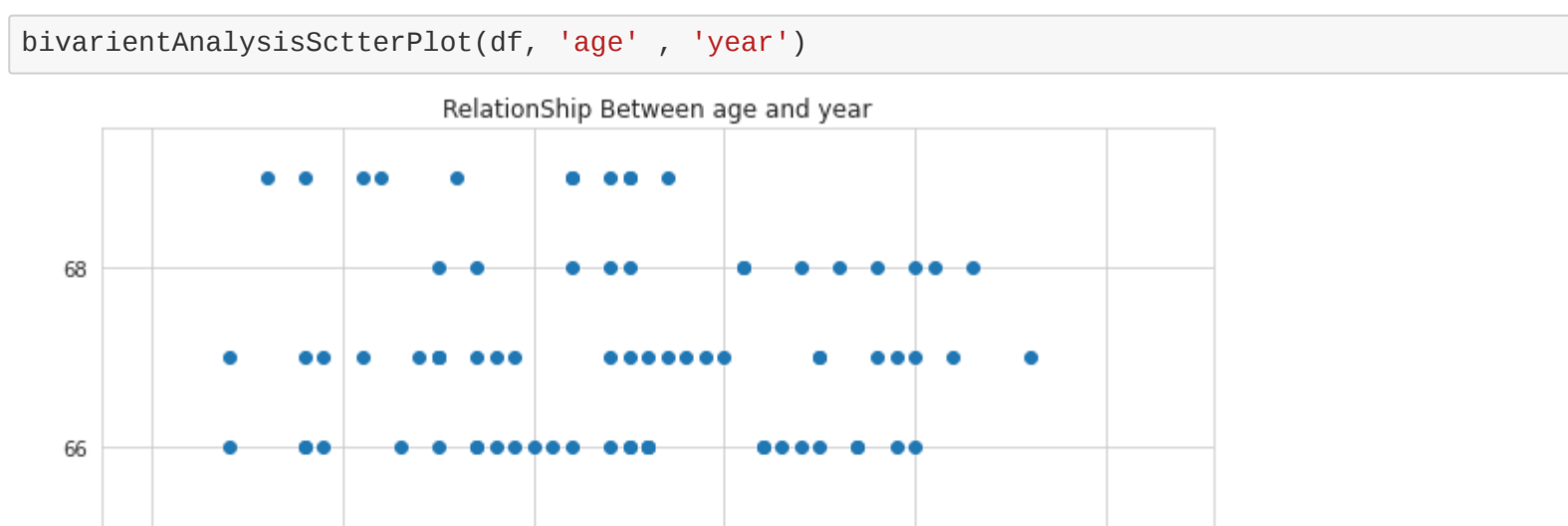
```
In [98]: bivariateAnalysisScatterPlot(df, 'age', 'year')
```



Box Plot is important for univariate analysis, this Plot Shows more no to people survived in 42- 52 age (25% of people survived more than five year in this range)

## Violin Plot

```
In [103]: fig, axes = plt.subplots(1, 3, figsize=(15, 5))
for idx, feature in enumerate(list(df.columns)[1:-1]):
    sns.violinplot(x='status', y=feature,
data=df, ax=axes[idx])
plt.show()
```



Observe the 3rd subplot See More no of people are survived only affected less than 5 positive lymph node, As we saw the orange plot the patient whose not survived has affect more positive lymph node, as per understanding if the a patient affected more lymph positive node chance of survival is more than five is too low