

Haberman's Survival Dataset

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.

```
In [192]: import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
import numpy as np
from sklearn.preprocessing import StandardScaler

# Load haberman.csv into a pandas DataFrame
hm = pd.read_csv('haberman.csv')
hm.head(10)

Out[192]:
```

age	year	nstatus	status
0	30	66	1
1	30	62	3
2	30	65	0
3	31	59	2
4	31	65	4
5	33	58	10
6	33	60	0
7	34	59	0
8	34	66	9
9	34	58	30

```
In [26]: # (q) How many data points and features?
print(hm.shape)
(306, 4)

In [27]: # (q) What are the column names in our dataset?
print(hm.columns)
Index(['age', 'year', 'nodes', 'status'], dtype='object')

In [28]: # (q) How many data points for each class are present?
hm['status'].value_counts()
#haberman is an unbalanced dataset

Out[28]:
```

status	count
1	225
2	81

name: status, dtype: int64

Attribute Information

There are 306 examples and 4 features in this dataset.

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

Four features have age, year and nodes are independent features while status as dependent feature.

Objective

In this dataset, using features age, year and nodes, we have to classify results in two classes-land 2

```
In [30]: # 2-D Scatter plot nodes vs age
hm.plot(kind='scatter', x='age', y='nodes')
plt.show()
```

Observation

This scatter plot does not give any relevant information regarding our dataset. \ So we will go for different scatter plots with color coding. We color the points by their class-label.

```
In [44]: # 2-D Scatter plot with color-coding for each class.
# How many combinations exist? 3C2 = 3.
# Year vs age
sns.set_style('whitegrid')
sns.FacetGrid(hm, hue='status', height=6)\
    .map(plt.scatter, 'age', 'nodes')\
    .add_legend()
plt.title('Year vs Age')
plt.show()
```

```
In [45]: # 2-D Scatter plot with color-coding for Nodes vs Age.
sns.set_style('whitegrid')
sns.FacetGrid(hm, hue='status', height=6)\
    .map(plt.scatter, 'age', 'nodes')\
    .add_legend()
plt.title('Nodes vs Age')
plt.show()
```

```
In [46]: # 2-D Scatter plot with color-coding for Nodes vs Year.
sns.set_style('whitegrid')
sns.FacetGrid(hm, hue='status', height=6)\
    .map(plt.scatter, 'year', 'nodes')\
    .add_legend()
plt.title('Nodes vs Year')
plt.show()
```

Observation

1. orange and blue data points cannot be easily separated in each case. 2. Separating class 1 from class 2 is much harder as they have considerable overlap.

```
In [49]: # pairwise scatter plot: Pair-Plot
plt.close();
sns.set_style('whitegrid');
sns.pairplot(hm, hue='status', height=3, vars=['age', 'year', 'nodes'])
plt.show()
# NOTE: the diagonal elements are PDFs for each feature. PDFs are explained below.
```

```
In [55]: # 1-D scatter plot of age
import numpy as np
one = hm.loc[hm['status'] == 1]
two = hm.loc[hm['status'] == 2]

# Print iris dataset [petal, length]
plt.plot(one['age'], np.zeros_like(one['age']), 'o')
plt.plot(two['age'], np.zeros_like(two['age']), 'o')
plt.legend('12')
plt.title('Age in Class 1 and 2')
plt.xlabel('Age')
plt.show()
```

```
In [58]: # 1-D scatter plot of year
import numpy as np
one = hm.loc[hm['status'] == 1]
two = hm.loc[hm['status'] == 2]

# Print iris dataset [petal, length]
plt.plot(one['year'], np.zeros_like(one['year']), 'o')
plt.plot(two['year'], np.zeros_like(two['year']), 'o')
plt.legend('12')
plt.title('Year in Class 1 and 2')
plt.xlabel('Year')
plt.show()
```

```
In [59]: # 1-D scatter plot of nodes
import numpy as np
one = hm.loc[hm['status'] == 1]
two = hm.loc[hm['status'] == 2]

# Print iris dataset [petal, length]
plt.plot(one['nodes'], np.zeros_like(one['nodes']), 'o')
plt.plot(two['nodes'], np.zeros_like(two['nodes']), 'o')
plt.legend('12')
plt.title('Nodes in Class 1 and 2')
plt.xlabel('Nodes')
plt.show()
```

Observation

From the above Pair-plot, we are not able to separate any of the data points. \ So we cannot classify whether the given observation is of class 1 or class 2. \ 1-D plots are also not giving any relevant information. So we will go for the histogram plots of features.

```
In [61]: #Histogram and PDF of age of Class 1 and Class 2
sns.FacetGrid(hm, hue='status', height=5)\
    .map(sns.distplot, 'age')\
    .add_legend()
plt.title('Histogram and PDF of Age of Class 1 and 2')
plt.show()
```

```
In [63]: #Histogram and PDF of age of Class 1 and Class 2
sns.FacetGrid(hm, hue='status', height=5)\
    .map(sns.distplot, 'year')\
    .add_legend()
plt.title('Histogram and PDF of Age of Class 1 and 2')
plt.show()
```

```
In [64]: #Histogram and PDF of age of Class 1 and Class 2
sns.FacetGrid(hm, hue='status', height=5)\
    .map(sns.distplot, 'nodes')\
    .add_legend()
plt.title('Histogram and PDF of Age of Class 1 and 2')
plt.show()
```

Observation

From the above histograms of age and year, we can see that PDFs of both the classes are overlapped. So we cannot consider any of the two features for classification. But in case of histogram plot of nodes, we can see that even though there is a little bit overlap between PDFs of classes 1 and 2 but there is a possibility of making a classifier based on the feature when compared to other two features. So we can make use of this feature only.

```
In [73]: #CDF and PDF plots of nodes of class 1
counts, bin_edges = np.histogram(one['nodes'], bins=10, density = True)
pdf = counts/sum(counts)
print(bin_edges)
print(pdf)
cdf = np.cumsum(pdf)
plt.plot(bin_edges[1:], pdf, label='pdf of class 1')
plt.plot(bin_edges[1:], cdf, label='cdf of class 1')
plt.legend()
plt.show()
```

```
In [73]: #CDF and PDF plots of nodes of class 2
counts, bin_edges = np.histogram(two['nodes'], bins=10, density = True)
pdf = counts/sum(counts)
print(bin_edges)
print(pdf)
cdf = np.cumsum(pdf)
plt.plot(bin_edges[1:], pdf, label='pdf of class 2')
plt.plot(bin_edges[1:], cdf, label='cdf of class 2')
plt.legend()
plt.show()
```

```
In [76]: #CDF plots of nodes of Class 1 and Class 2
counts, bin_edges=np.histogram(one['nodes'], bins=10, density=True)
pdf=counts/sum(counts)
cdf=np.cumsum(pdf)
concd=plt.plot(bin_edges[1:],cdf,label='CDF Class 1')
counts,bin_edges=np.histogram(two['nodes'], bins=10, density=True)
pdf=counts/sum(counts)
cdf=np.cumsum(pdf)
twocdf=plt.plot(bin_edges[1:],cdf,label='CDF Class 2')
plt.xlabel('Nodes')
plt.ylabel('Probability')
plt.title('CDF of Nodes of Classes 1 and 2')
plt.legend()
plt.show()
```

Observation

From the above CDF curves of nodes we get to know that about 84% of Class 1 people have nodes less than 5 and about 45% of Class 2 people have nodes less than 5. \ So we have got acclissar boundary at nodes=5. Now we will perform statistical operations on feature.

Means:

```
In [80]: print("Mean of age")
print("Class 1=",np.mean(one['age']))
print("Class 2=",np.mean(two['age']))

Mean of age
Class 1= 52.01777777777778
Class 2= 35.67501234567901

In [81]: print("Mean of year")
print("Class 1=",np.mean(one['year']))
print("Class 2=",np.mean(two['year']))

Mean of year
Class 1= 61.80222222222222
Class 2= 62.8716049382716

In [82]: print("Mean of nodes")
print("Class 1=",np.mean(one['nodes']))
print("Class 2=",np.mean(two['nodes']))

Mean of nodes
Class 1= 2.7911111111111113
Class 2= 7.45578012345679
```

Std-dev:

```
In [85]: print("Standard Deviation of age")
print("Class 1=",np.std(one['age']))
print("Class 2=",np.std(two['age']))

Standard Deviation of age
Class 1= 10.8876547510301
Class 2= 11.1018131301113

In [86]: print("Standard Deviation of year")
print("Class 1=",np.std(one['year']))
print("Class 2=",np.std(two['year']))

Standard Deviation of year
Class 1= 3.13745314601396
Class 2= 3.321423625020783

In [84]: print("Standard Deviation of nodes")
print("Class 1=",np.std(one['nodes']))
print("Class 2=",np.std(two['nodes']))

Standard Deviation of nodes
Class 1= 0.0
Class 2= 9.12874676761632
```

Medians:

```
In [86]: print("Median of age")
print("Class 1=",np.median(one['age']))
print("Class 2=",np.median(two['age']))

Median of age
Class 1= 54.0
Class 2= 55.0

In [87]: print("Median of year")
print("Class 1=",np.median(one['year']))
print("Class 2=",np.median(two['year']))

Median of year
Class 1= 61.0
Class 2= 63.0

In [88]: print("Median of nodes")
print("Class 1=",np.median(one['nodes']))
print("Class 2=",np.median(two['nodes']))

Median of nodes
Class 1= 0.0
Class 2= 4.0
```

Quantile

```
In [89]: print("Unquantiles:")
print(np.percentile(one['age'], np.arange(0,100,25)))
print(np.percentile(two['age'], np.arange(0,100,25)))

Quantiles:
[30. 45. 50. 60.]
[34. 46. 53. 61.]

In [90]: print("Unquantiles:")
print(np.percentile(one['year'], np.arange(0,100,25)))
print(np.percentile(two['year'], np.arange(0,100,25)))

Quantiles:
[58. 65. 63. 66.]
[58. 59. 63. 65.]

In [91]: print("Unquantiles:")
print(np.percentile(one['nodes'], np.arange(0,100,25)))
print(np.percentile(two['nodes'], np.arange(0,100,25)))

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

MAD

```
In [93]: print("Median Absolute Deviation")
print(robust.mad(one['age']))
print(robust.mad(two['age']))

Median Absolute Deviation
13.3438194655017
31.86081748044816

In [94]: print("Median Absolute Deviation")
print(robust.mad(one['year']))
print(robust.mad(two['year']))

Median Absolute Deviation
4.4478665516806
4.4478665516806

In [95]: print("Median Absolute Deviation")
print(robust.mad(one['nodes']))
print(robust.mad(two['nodes']))

Median Absolute Deviation
0.0
5.93048874022408
```

Observation

These statistical operations give us a lot of information. Means, Medians, Standard Deviations and Quantiles of features, age and year are almost same. So, we cannot differentiate between them. But in case of nodes, they show significant differences. So again it tells we can use nodes to make classifier boundary. \ To differentiate, we will plot boxplot.

```
In [96]: #Box Plot of age of Class 1 and Class 2
sns.boxplot(x='status', y='age', data=hm)
plt.title('Box Plot of age of Class 1 and Class 2')
plt.show()
```

```
In [104]: #Box Plot of year of Class 1 and Class 2
sns.boxplot(x='status', y='year', data=hm)
plt.title('Box Plot of age of Class 1 and Class 2')
plt.show()
```

```
In [106]: #Box Plot of nodes of Class 1 and Class 2
sns.boxplot(x='status', y='nodes', data=hm)
plt.title('Box Plot of age of Class 1 and Class 2')
plt.show()
```

```
In [107]: #Violin Plot of age of Class 1 and Class 2
sns.violinplot(x='status', y='age', data=hm)
plt.title('Violin Plot of age of Class 1 and Class 2')
plt.show()
```

```
In [108]: #Violin Plot of year of Class 1 and Class 2
sns.violinplot(x='status', y='year', data=hm)
plt.title('Violin Plot of age of Class 1 and Class 2')
plt.show()
```

Obsevation

Again from the boxplots and violin plots, it is clear that features age and year do not help in classifying as their quantiles do not show very much difference. But from the boxplot of feature nodes, we can make certain inferences as the 75th percentile of class 1 is less than 50th percentile of class 2 which is about 5. This means that 75 percent of nodes in class 1 are less than or equal to 5 while in class 2 we have only about 50 percent of nodes which are less than or equal to 5.

Conclusion