	 age: Petiont's Age year: Operation Year(1958-1970) nodes: infected Axillary nodes detected in petiont. status: Whether petiont has Survived or Not Survived #Replacing Variables for better understanding df. status[df['status']==1]='Survived' df. status[df['status']==2]='NotSurvived' df. head() #overview of DataSet age year nodes status 0 30 62 3 Survived
]:	1 30 65 0 Survived 2 31 59 2 Survived 3 31 65 4 Survived 4 33 58 10 Survived #Description of DataSet print(df.describe()) print("\nshape ", df.shape, "\n") print(df.status.value_counts(), "\n\n") age
	Count 305.000000 305.000000 305.000000 305.000000 mean 52.531148 62.849180 4.036066 std 10.744024 3.254078 7.199370 min 30.000000 58.00000 0.0000000 258.000000 0.0000000 65.000000 63.000000 63.000000 66.000000 66.000000 66.000000 66.000000 52.000000 max 83.000000 69.00000 52.000000 stape (305, 4) Survived 224 NotSurvived 81 Name: status, dtype: int64
	By observing the records in given DataSet The average age of petionts is 52 years. The Youngest Petiont is of 30 years old,whereas Eldest is 83 years old. The survival rate of petionts is about 73% (224/305*100=73.44%) According to above observation the diffrence between above classes is found to be higher. Hence,our DataSet is imbalanced. Univariate Analysis
	Univariate Analysis PDF If a Function follows continuous Destribution, The PDF (probability Density function) is a probability that a random variable X will take value exactly equal to x(a particular value) In simpler words, PDF is a percetage of data at that point. sns.FacetGrid(df, hue='status', size=5).map(sns.distplot, 'year').add_legend() plt.ylabel('P(X)=x') plt.title('PDF of year') plt.show() PDF of year
	0.08 0.08 0.06 0.04 0.02 status Survived NotSurvived
	Difficult to conclude the plot due to overlap. sns.FacetGrid(df, hue='status', size=5).map(sns.distplot, 'age').add_legend() plt.title('PDF of age') plt.ylabel('P(X)=x') print("\n\n") sns.FacetGrid(df, hue='status', size=5).map(sns.distplot, 'nodes').add_legend() plt.ylabel('P(X)=x') print("\n\n") sns.FacetGrid(df, hue='status', size=5).map(sns.distplot, 'nodes').add_legend() plt.ylabel('P(X)=x') plt.title('PDF of nodes') plt.title('PDF of nodes') plt.show()
	0.030 0.025 0.020 0.015 PDF of age status Survived NotSurvived
	0.005 0.000 20 30 40 50 60 70 80 90 PDF of nodes 0.5
	0.4 0.2 0.1 0.0 -10 0 10 20 30 40 50 60 nodes
	 From above two plots it can be observed 1. people having age less than or equal to 39 have higher chances of survival ,wheras it gets inversed with people aged greater than 39.(statement is unclear) 1. lesser the cancer spread across nodes higher the chances survival, the statement has found out to be correct(As most of the plot have verified or concluded our statement.) 1. Chances are sligtly higher that the person will not survive if cancer has spread to more than 4 nodes. As, more the conclusion come out of Data better it iswe will try getting more conclusions by plotting few more plots. CDF The Cumulative Destribution Function is the probability that a random variable X will take value lesser than equals to x(perticular value). In simpler words, CDF tells the percetage of data lesser than equal to that point(x).
	<pre>counts, bins=np.histogram(df['age'], bins=18, density=True) pdf=counts/(sum(counts)) cdf=np.cumsum(pdf) plt.plot(bins[1:], pdf, label='pdf') plt.plot(bins[1:], cdf, label='cdf') plt.xlabel("age") plt.ylabel("PDF: P(X)=x\nCDF: P(X)<=x") plt.legend(bbox_to_anchor=(1.3, 0.6)) plt.title("PDF & CDF of age") plt.show() print(f"Survived:\n CDF:{cdf} \n BinEdges:{bins} \n\n\n\n\n\n\n")</pre>
	PDF & CDF of age 10 08 K No
	Survived: CDF:[0.01311475 0.04918033 0.10819672 0.1704918 0.25901639 0.34754098 0.44262295 0.5442623 0.64262295 0.72786885 0.8 0.87540984 0.91803279 0.95737705 0.98360656 0.99344262 0.99672131 1
	Around 1% of petionts have age lesser than 30 years. Around 44% of petionts have age lesser than 50 years. Around 80% of petionts have age lesser than 71 years. Around 99% of petionts have age lesser than 80 years. counts, bins=np.histogram(df['year'], bins=10, density=True) pdf=counts/(sum(counts)) cdf=np.cumsum(pdf)
	plt.plot(bins[1:], pdf, label='pdf') plt.plot(bins[1:], cdf, label='cdf') plt.xlabel("year") plt.ylabel("PDF: P(X)=x\nCDF: P(X)<=x") plt.legen(bbox_to_anchor=(1.3, 0.6)) plt.title("PDF & CDF of year") plt.show() print(f"Survived:\n CDF:{cdf} \n BinEdges:{bins} \n\n\n\n\n\n") PDF & CDF of year 10 08
	Survived: CDF:[0.20655738 0.29836066 0.38360656 0.45901639 0.55737705 0.6557377 0.74754098 0.83934426 0.92131148 1. BinEdges:[58. 59.1 60.2 61.3 62.4 63.5 64.6 65.7 66.8 67.9 69.]
	Around 55% of petionts had been operated year lesser than or equal to 63. Due to the limited number of operation years in DataSet,We are unable to draw many conclusions. Survived=df['nodes'][df['status']=='Survived'] NotSurvived=df['nodes'][df['status']=='NotSurvived'] counts, bins=np.histogram(Survived, bins=18, density=True) pdf=counts/(sum(counts)) cdf=np.cumsum(pdf)
	<pre>plt.plot(bins[1:], pdf, label='S_pdf') plt.plot(bins[1:], cdf, label='S_cdf') plt.xlabel("Nodes") plt.ylabel("PDF: P(X)=x\nCDF: P(X)<=x") plt.legend(bbox_to_anchor=(1.3, 0.6)) plt.title("The PDF and CDF of nodes for Survived") plt.show() print(f"Survived:\n CDF:{cdf} \n BinEdges:{bins} \n\n\n\n\n\n") print("\n\n The PDF and CDF of nodes for NotSurvived") counts1, bins1=np.histogram(NotSurvived, bins=18, density=True) NSpdf=counts1/(sum(counts1))</pre>
	NScdf=np.cumsum(NSpdf) plt.plot(bins1[1:],NSpdf,label='NS_pdf') plt.plot(bins1[1:],NScdf,label='NS_cdf') plt.xlabel("Nodes") plt.ylabel("PDF: P(X)=x\nCDF: P(X)<=x") plt.legend(bbox_to_anchor=(1.3, 0.6)) plt.title("The PDF and CDF of nodes for NotSurvived") plt.show() print(f"NotSurvived:\n CDF:{NScdf} \n BinEdges:{bins1}\n") The PDF and CDF of nodes for Survived
	Survived: CDF:[0.73214286 0.84375
	0.99553571 0.99553571 0.99553571 0.99553571 1.] BinEdges:[0.
	10 08 08 08 06 02 00 00 00 00 00 00 00 00 00 00 00 00
]:	CDF:[0.39506173 0.56790123 0.65432099 0.75308642 0.82716049 0.86419753 0.90123457 0.96296296 0.97530864 0.97530864 0.97530864 0.97530864 0.98765432 0.98765432 0.98765432 0.98765432 1. BinEdges:[0. 2.8888889 5.77777778 8.66666667 11.55555556 14.44444444 17.33333333 20.22222222 23.11111111 26. 28.8888889 31.77777778 34.66666667 37.55555556 40.44444444 43.33333333 46.22222222 49.11111111 52.] #Values are evaluated by slightly changing bin counts. Tabular representation: Nodes Survived NotSurvived Difference
	2 73% 39% 34% 5 84% 56% 28% 10 92% 71% 21% 15 95% 85% 10% 30 99% 97% 3% Observation Observing the difference, It has gotten clearer that The petionts chances of Survival decline with higher number of infected nodes
	sns.boxplot(x='status',y='year',data=df) plt.title('Boxplot of status-year') plt.grid plt.show() Boxplot with status-year
)]:	sns.boxplot(x='status', y='age', data=df) plt.title('BoxPlot of age-status') plt.grid plt.show() BoxPlot with age-status
	80 80 90 Survived Status NotSurvived Major overlap can be observed.
	sns.boxplot(x='status',y='nodes',data=df) plt.title('Boxplot of nodes-status') plt.show() Boxplot with nodes-status 50 40 82 20 40 40 40 40 40 40 40 40 4
	#(df['age'][df['status']=='NotSurvived']).min() For Survived, • The box plot with survived has values below 50th percentile all lying on 0,and therefore most of the petionts survived have 0 infected nodes. • 75th percentile of petionts survived have about 3 or 4 infected nodes.
	 The points beyond whiskers are less likely to be outliers because there are more than normal number of them.But the point at 46(i.e,Petiont with 46 infected nodes) is an extreme outlier,Because person generally ha about 20-30 axillary nodes,person having 46 nodes is not possible. For NotSurvived, 25th percentile of petionts not survived had about 0 or 1 infected nodes. About 50th percentile of petionts had infected nodes lesser than or equal to 4(<=: because the median lies on 4.5) and 75th percentile of petionts have lesser than equal to 11 infected nodes. The maximum number of infected nodes in petiont are 24.(excluding the extreme outliers). VIOLIN PLOTS sns. violinplot (data=df, x='year', y='nodes', size=30) plt. title('Violin plot of year-nodes') plt. show()
	sns.violinplot(data=df, x='year', y='age', size=30) plt.title('Violin plot of year-age') plt.show() sns.violinplot(data=df, x='age', y='status', size=30) plt.title('Violin plot of age-status') plt.show() Violin plot of year-nodes Violin plot of year-nodes
	Violin plot of year-age
	70 40 30 20 58 59 60 61 62 63 64 65 66 67 68 69 year Violin plot of age-status
·]:	Survived NotSurvived 20 30 40 50 50 60 70 80 90 #df['nodes'][df['status']=='NotSurvived'][df['nodes']>=24]
	# We are undable to make sense of data because major overlap can be observed. #So, We will try drawing conclusions through overall plot. On plot of nodes-year: We can observe the densities of plots are desributed mostly on points lesser than 11 and even more on 0 and 1, It shows most of the petionts have lesser than 10 infected nodes and many of them have about 0 or 1 infected nodes. On plot of age-year: We can observe that the densities of plots are destributed mostly on points from 45 to 60, It shows most of petionts have age around 45 to 60. On plot of age-status: We can observe that around 50 percent of people were about 43 to 62 year old.
	Bi-Variete Analysis SCATTER PLOT sns.set_style('whitegrid') sns.FacetGrid(df,hue='status',size=5).map(plt.scatter,'age','nodes').add_legend().set(title='Scatter Plot of nodes-age') plt.show() print('\n') sns.FacetGrid(df,hue='status',size=5).map(plt.scatter,'age','year').add_legend().set(title='Scatter plot of year-age') plt.show() print('\n') sns.FacetGrid(df,hue='status',size=5).map(plt.scatter,'year','nodes').add_legend().set(title='Scatter plot of nodes-year') plt.show()
	10 30 40 50 80 70 80 Scatter plot of year-age
	66 62 80 80 80 80 80 80 80 80 80 80 80 80 80
	Scatter plot of nodes-year 50 40 30 Status Survived NotSurvived
	Plotwise observation Plot of nodes-age and year-age: • Irrespective of nodes and year people having age lesser than 40 have higher chances of survival.
	Plot of nodes-age and nodes-year: • Many of petionts have lesser than 10 infected nodes & Most of them have 0 infected nodes. We can observe that for observation of age, the age-year plot shows clearer representation. Whereas,for observation of number of nodes the nodes-age plot shows better representation. PAIR PLOT #Pairplot sns.set_style("whitegrid"); sns.pairplot(df, hue='status', size=5); plt.show()
	<pre>sns.pairplot(df,hue='status',size=5);</pre>
	40 30 68 68 68 68 68 68 68 68 68 6
	64 62 62 63 64 65 65 65 65 65 65 65 65 65 65 65 65 65
	50
	The scatter plot with age and nodes shows: 1.lesser the cancer spread across nodes, higher the chances of survival. 2.Plot of nodes-age shows better representation.
	Overall Conclusion through Analysis: • Through Description: 1. The Youngest Petiont is of 30 years old, whereas Eldest is 83 years old. 2. The survival rate of petionts is about 73%. • Through Uni-Variate & Bi-Variate Analysis: 1. petiont has higher chances of survival if petiont is lesser than 40 year old or petiont has lesser age. 2. petiont has lesser infected nodes, even better if petiont has lesser than 4 infected nodes. 3. Most of the surviving petionts have 0 infected nodes. 4. Most of the petionts (irrespective of status) have lesser than 11 infected nodes. 5. Most of petionts(around 50%) are 43 to 63 years old
7]:	<pre>%%html <style> table {float:left} </style></pre>