Exploratory Data Analysis of Haberman's **DataSet** In [109]: # Install the PyDrive wrapper & import libraries. # This only needs to be done once in a notebook. !pip install -U -q PyDrive from pydrive.auth import GoogleAuth from pydrive.drive import GoogleDrive from google.colab import auth from oauth2client.client import GoogleCredentials # Authenticate and create the PyDrive client. # This only needs to be done once in a notebook. auth.authenticate\_user() gauth = GoogleAuth() gauth.credentials = GoogleCredentials.get application default() drive = GoogleDrive(gauth) # Create & upload a text file. uploaded = drive.CreateFile({'title': 'Sample file.txt'}) uploaded.SetContentString('Sample upload file content') uploaded.Upload() print('Uploaded file with ID {}'.format(uploaded.get('id'))) Uploaded file with ID 1zi619Ah7wg BIKGm1tf1cEMsfvCxsiJJ In [110]: # Download a file based on its file ID. file id = '1YJzytNcpvROyRBGC lRt nUezp0dczae' downloaded = drive.CreateFile({'id': file\_id}) 'Downloaded content "{}"'.format(downloaded.GetContentString()) # https://drive.google.com/open?id=1YJzytNcpvROyRBGC lRt nUezpOdczae Out[110]: 'Downloaded content "age, operation\_year, axil\_nodes, status\n30,64,1,1 \n30,62,3,1\n30,65,0,1\n31,59,2,1\n31,65,4,1\n33,58,10,1\n33,60,0,1\n  $34,59,0,2 \times 34,66,9,2 \times 34,58,30,1 \times 34,60,1,1 \times 34,61,10,1 \times 34,67,7,1 \times 34,59,0,2 \times 34$  $4,60,0,1 \times 35,64,13,1 \times 35,63,0,1 \times 36,60,1,1 \times 36,69,0,1 \times 37,60,0,1 \times 37,6$  $63,0,1 \times 37,58,0,1 \times 37,59,6,1 \times 37,60,15,1 \times 37,63,0,1 \times 38,69,21,2 \times 38,5$ 9,2,1\n38,60,0,1\n38,60,0,1\n38,62,3,1\n38,64,1,1\n38,66,0,1\n38,66,1  $1,1\n38,60,1,1\n38,67,5,1\n39,66,0,2\n39,63,0,1\n39,67,0,1\n39,58,0,1$ \n39,59,2,1\n39,63,4,1\n40,58,2,1\n40,58,0,1\n40,65,0,1\n41,60,23,2\n 41,64,0,2 n41,67,0,2 n41,58,0,1 n41,59,8,1 n41,59,0,1 n41,64,0,1 n41,69,8,1\n41,65,0,1\n41,65,0,1\n42,69,1,2\n42,59,0,2\n42,58,0,1\n42,60,  $1,1 \cdot 1,1 \cdot 1,1$  $2,2\n43,59,2,2\n43,64,0,2\n43,63,14,1\n43,64,2,1\n43,64,3,$  $1\n43,60,0,1\n43,63,2,1\n43,65,0,1\n43,66,4,1\n44,64,6,2\n44,58,9,2\n$ 44,63,19,2\n44,61,0,1\n44,63,1,1\n44,61,0,1\n44,67,16,1\n45,65,6,2\n4 5,66,0,2\n45,67,1,2\n45,60,0,1\n45,67,0,1\n45,59,14,1\n45,64,0,1\n45, 68,0,1\n45,67,1,1\n46,58,2,2\n46,69,3,2\n46,62,5,2\n46,65,20,2\n46,6  $2,0,1 \land 46,58,3,1 \land 63,0,1 \land 63,23,2 \land 62,0,2 \land 65,0,2 \land 61,$  $0,1\\n47,63,6,1\\n47,66,0,1\\n47,67,0,1\\n47,58,3,1\\n47,60,4,1\\n47,68,4,1$ \n47,66,12,1\n48,58,11,2\n48,58,11,2\n48,67,7,2\n48,61,8,1\n48,62,2,1  $\n48,64,0,1\n48,66,0,1\n49,63,0,2\n49,64,10,2\n49,61,1,1\n49,62,0,1\n$  $49,66,0,1 \\ \\ 1,1 \\ 1,$ 63,13,2\n50,64,0,2\n50,59,0,1\n50,61,6,1\n50,61,0,1\n50,63,1,1\n50,5  $8,1,1 \in \{0,0\}$  $3,2\n51,59,3,2\n51,64,7,1\n51,59,1,1\n51,65,0,1\n51,66,1,1\n52,69,3,2$  $\n52, 59, 2, 2 \n52, 62, 3, 2 \n52, 66, 4, 2 \n52, 61, 0, 1 \n52, 63, 4, 1 \n52, 69, 0, 1 \n52, 63, 4, 2 \n52, 61, 0, 2 2 \n52$  $2,60,4,1 \n52,60,5,1 \n52,62,0,1 \n52,62,1,1 \n52,64,0,1 \n52,65,0,1 \n52,6$  $8,0,1 \n53,58,4,2 \n53,65,1,2 \n53,59,3,2 \n53,60,9,2 \n53,63,24,2 \n53,65,$  $12,2\n53,58,1,1\n53,60,1,1\n53,60,2,1\n53,61,1,1\n53,63,0,1\n54,60,1$  $1,2\n54,65,23,2\n54,65,5,2\n54,68,7,2\n54,59,7,1\n54,60,3,1\n54,66,0,$ 1\n54,67,46,1\n54,62,0,1\n54,69,7,1\n54,63,19,1\n54,58,1,1\n54,62,0,1 \n55,63,6,2\n55,68,15,2\n55,58,1,1\n55,58,0,1\n55,58,1,1\n55,66,18,1 \n55,66,0,1\n55,69,3,1\n55,69,22,1\n55,67,1,1\n56,65,9,2\n56,66,3,2\n 56,60,0,1\n56,66,2,1\n56,66,1,1\n56,67,0,1\n56,60,0,1\n57,61,5,2\n57, 62,14,2\n57,64,1,2\n57,64,9,1\n57,69,0,1\n57,61,0,1\n57,62,0,1\n57,6  $3,0,1 \n57,64,0,1 \n57,64,0,1 \n57,67,0,1 \n58,59,0,1 \n58,60,3,1 \n58,61,$  $1,1\n58,67,0,1\n58,58,0,1\n58,58,3,1\n58,61,2,1\n59,62,35,2\n59,60,0,$  $1\n59,63,0,1\n59,64,1,1\n59,64,4,1\n59,64,0,1\n59,64,7,1\n59,67,3,1\n$ 60,59,17,2\n60,65,0,2\n60,61,1,1\n60,67,2,1\n60,61,25,1\n60,64,0,1\n6 1,62,5,2\n61,65,0,2\n61,68,1,2\n61,59,0,1\n61,59,0,1\n61,64,0,1\n61,6  $5,8,1 \n61,68,0,1 \n61,59,0,1 \n62,59,13,2 \n62,58,0,2 \n62,65,19,2 \n62,6$ 2,6,1\n62,66,0,1\n62,66,0,1\n62,58,0,1\n63,60,1,2\n63,61,0,1\n63,62,  $0,1 \n63,63,0,1 \n63,63,0,1 \n63,66,0,1 \n63,61,9,1 \n63,61,28,1 \n64,58,0,$ 1\n64,65,22,1\n64,66,0,1\n64,61,0,1\n64,68,0,1\n65,58,0,2\n65,61,2,2 \n65,62,22,2\n65,66,15,2\n65,58,0,1\n65,64,0,1\n65,67,0,1\n65,59,2,1 \n65,64,0,1\n65,67,1,1\n66,58,0,2\n66,61,13,2\n66,58,0,1\n66,58,1,1\n 66,68,0,1\n67,64,8,2\n67,63,1,2\n67,66,0,1\n67,66,0,1\n67,61,0,1\n67, 65,0,1\n68,67,0,1\n68,68,0,1\n69,67,8,2\n69,60,0,1\n69,65,0,1\n69,66,  $0,1\n70,58,0,2\n70,58,4,2\n70,66,14,1\n70,67,0,1\n70,68,0,1\n70,59,8,$  $1 \times 70,63,0,1 \times 71,68,2,1 \times 72,63,0,2 \times 72,58,0,1 \times 72,64,0,1 \times 72,67,3,1 \times 72,63,0,1 \times 72,0,1 \times$  $73,62,0,1 \times 73,68,0,1 \times 74,65,3,2 \times 74,63,0,1 \times 75,62,1,1 \times 76,67,0,1 \times 77,$ 65,3,1\n78,65,1,2\n83,58,2,2\n"' In [0]: | # Ignoring the warnings import warnings warnings.filterwarnings('ignore') In [0]: # Getting the content from the csv file downloaded.GetContentFile('haberman.csv') In [113]: # Checking whether csv file is loaded or not adc.json haberman.csv sample data In [0]: # Importing the required libraries needed for the analysis import numpy as py import pandas as pd import matplotlib.pyplot as plt import seaborn as sb sb.set() In [115]: # Using pandas creating a data frame df = pd.read csv('haberman.csv') df.head() Out[115]: operation\_year axil\_nodes status 0 30 1 64 1 1 30 62 3 1 0 2 30 65 1 3 31 59 2 1 1 31 65 In [116]: # Colums present in the data set df.columns Out[116]: Index(['age', 'operation\_year', 'axil\_nodes', 'status'], dtype='objec t') In [117]: # Finding the number of datapoins, feaatures df.shape Out[117]: (306, 4) **Obervation:**  306 is the row count which indiactes the number of patients and 4 is the column count which indicates the variable count # Count of similar kind of data in the patient status In [118]: df['status'].value\_counts() Out[118]: 225 81 Name: status, dtype: int64 Obervation: 225 and 81 shows the count of patients who have survived more than 5 years and less than 5 years after the operation respectively. **Bivariate Analysis** Scatterplot In [140]: #plotting the scatter plot df.plot(kind = 'scatter', x = 'axil\_nodes', y = 'age') plt.title('scatter plot - age vs axil\_nodes',color = 'y',size = '13') plt.xlabel('axil\_nodes', size = 13, color = 'b') plt.ylabel('age', size = 13, color = 'b') plt.show() 'c' argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its 1 ength matches with 'x' & 'y'. Please use a 2-D array with a single r ow if you really want to specify the same RGB or RGBA value for all p oints. No handles with labels found to put in legend. scatter plot - age vs axil\_nodes 80 70 60 50 40 30 axil\_nodes Obervation: From the above scatter plots we can infer that most the patients dont have axil\_nodes (i.e axil\_nodes =0 )and very few patients have axil\_nodes above 30 Swarm plot #plotting the swarm plot where data about axil nodes, age and status can In [141]: be represented. sb.swarmplot(x = 'axil nodes', y = 'age', data = df, hue = 'status', size plt.title('swarm plot - age vs axil nodes',color = 'y',size = '13') plt.xlabel('axil nodes', size = 13, color = 'b') plt.ylabel('age', size = 13, color = 'b') plt.show() swarm plot - age vs axil\_nodes 1 80 70 50 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 28 30 35 46 52 axil\_nodes observation: • From the above plot, we can observe more number of patients have survived more than 5 years have 0 axil nodes. • More number of patients have axil\_nodes < 5 and very less number of patients having more axil\_nodes and in the plot. U can observe as the number of axial nodes increase the density of patients decreases gradually and patients having axial\_nodes above 30 are very few Pair plots In [121]: # plotting the 2-dimensional pair plots between diffrent variable sb.set style('whitegrid') sb.pairplot(df, hue = 'status', vars = ['age', 'operation year', 'axil node s'], size = 3).add legend() plt.show() 70 60 50 40 30 68 operation year 64 62 50 40 30 axii 20 10 0 operation\_year **Observations:**  As observed in the swarm plot we can make same observations from the age vs axil\_node pairplot. we cannot make any observation on the operation\_year vs status from the pair plot since most of the datapoints were overlapped Analysis using statistical properties of data In [122]: #observing some statistical properties of the data df.describe() Out[122]: operation\_year axil nodes status age count 306.000000 306.000000 306.000000 306.000000 52.457516 62.852941 4.026144 1.264706 mean 3.249405 7.189654 0.441899 std 10.803452 30.000000 58.000000 0.000000 1.000000 min 60.000000 25% 44.000000 0.000000 1.000000 52.000000 63,000000 **50%** 1.000000 1.000000 4.000000 **75%** 60.750000 65.750000 2.000000 83.000000 2.000000 max 69.000000 52.000000 **Observations**  Opeartions were happened between 1958 to 1969 The patients are aged between 30 and 83 The mean of the status is 1.26 which little bit closer to 1 than 2. So we can say that more than half of the people lived for more than 5 years after the operation There are patients with 0 axil\_nodes #creating data frames for the data by splitting according to the status In [0]: df gt5 = df.loc[df['status'] == 1] df lt5 =df.loc[df['status'] == 2] In [125]: df gt5.describe() Out[125]: age operation\_year axil\_nodes status count 225.000000 225.000000 225.000000 225.0 mean 52.017778 62.862222 2.791111 1.0 11.012154 5.870318 3.222915 0.0 std min 30.000000 58.000000 0.000000 1.0 0.000000 25% 43.000000 60.000000 1.0 **50%** 52.000000 63.000000 0.000000 1.0 60.000000 **75**% 77.000000 69.000000 46.000000 max 1.0 Observations about the patients who lived more than 5 years: • The max age is 77 Highest number of axil\_nodes found is 46 count of patients is 225 and mean of axil\_nodes is 2.79 In [126]: df lt5.describe() Out[126]: age operation\_year axil\_nodes status count 81.000000 81.000000 81.000000 81.0 mean 53.679012 62.827160 7.456790 2.0 std 10.167137 3.342118 9.185654 0.0 min 34.000000 58.000000 0.000000 2.0 **25%** 46.000000 59.000000 1.000000 2.0 **50%** 53.000000 63.000000 4.000000 2.0 **75%** 61.000000 65.000000 11.000000 2.0 max 83.000000 69.000000 52.000000 2.0 **Observations**  The mean of axil\_nodes is 7.456 which is almost 2.5 times more than the axil\_nodes we found in the patients who lived more than 5 years. so we can say that the patient having more than 46 axil\_nodes died with in 5 years after the operation The patient having age more than 77 years have not survived more than 5 years after the operation **Univariate Analysis Histograms and PDF** In [143]: #Histogram plotting and the probability distribution function sb.FacetGrid( data = df, hue = 'status', size = 5) \ .map(sb.distplot , 'axil nodes') \ .add legend() plt.title('Histogram', color = 'y', size = '13') plt.xlabel('axil nodes', size = 13, color = 'b') plt.show() Histogram 0.5 0.4 0.3 status 1 2 0.2 0.0 -10 axil\_nodes In [128]: sb.FacetGrid( data = df, hue = 'status', size = 5) \ .map(sb.distplot ,'operation\_year').add\_legend() plt.title('Histogram', color = 'y', size = '13') plt.xlabel('operation\_year', size = 13, color = 'b') plt.show() Histogram 0.12 0.10 0.08 status 0.06 0.04 0.02 0.00 57.5 60.0 62.5 65.0 67.5 70.0 operation\_year In [129]: sb.FacetGrid( data = df, hue = 'status', size = 5) \ .map(sb.distplot , 'age').add\_legend() plt.title('Histogram', color = 'y', size = '13') plt.xlabel('age', size = 13, color = 'b') plt.show() Histogram 0.040 0.035 0.030 0.025 0.020 2 0.015 0.010 0.005 0.000 100 **Observations**  From the above PDFs age vs status graph and operation\_year vs status graph didnt provide much insight because the data is overlapping in both plots.so we cant make a proper analysis But from the axil\_nodes vs status graph we can say that the probability of patients having '0' axil\_nodes is high and the patients having axil\_nodes grater than 30 are very less CDF In [130]: ct, edge = py.histogram(df['axil nodes'], bins=10, density = **True**) pdf = ct/(sum(ct))print(pdf); print(edge) cdf = py.cumsum(pdf)plt.axis(x='axil nodes') plt.plot(edge[1:],pdf) plt.plot(edge[1:], cdf) plt.title('cumilative distribution function',color = 'y',size = 15) plt.xlabel('axil\_nodes',color = 'b',size = 13) plt.legend(['probability distribution function','cumilative distributio n function']) plt.show() [0.77124183 0.09803922 0.05882353 0.02614379 0.02941176 0.00653595 0.00326797 0. 0.00326797 0.00326797] 5.2 10.4 15.6 20.8 26. 31.2 36.4 41.6 46.8 52. ] cumilative distribution function 10 0.8 0.6 probability distribution function cumilative distribution function 0.4 02 0.0 10 20 40 50 30 axil\_nodes Observation • By observing the values we can say that the percentage of patients with '0' axil\_nodes is 77 and with 52 axil\_nodes is 0.3 we can observe that the patients who are having axil\_nodes greater than 30 are very few. U can clearly observe that the graph is closer to the axis as the count of axil\_nodes are increasing **Box plots and Strip plot** box=sb.boxplot(y='age', x='status', data=df, width=0.5, In [218]: palette="colorblind", hue = 'status') box.legend(bbox\_to\_anchor=(1.2,0.75)) #plt.legend() # adding stripplot to boxplot box=sb.stripplot(y='age', x='status', data=df, jitter=True, marker='o', alpha=0.9, color='black') plt.xlabel('status',color = 'b',size = '13') plt.ylabel('age',color = 'b',size = '13') plt.title('Boxplot and Strip plot', color = 'y', size = 15) plt.show() Boxplot and Strip plot 80 70 60 50 40 30 1 2 status **Observations**  From the stripp plot we can say that more patients have lived for more than 5 years • 50% of patients have age less than 53 years 75% of patients have age less than 62 years violin plots In [178]: #Plotting the violin plots plt.subplot(131) sb.violinplot(x='status', y='age', data=df, marker='o', alpha=0.4, hue ='status') plt.xlabel('status',color = 'b',size = '13') plt.ylabel('age',color = 'b',size = '13') plt.title('Status vs age',color = 'b',size = 13) plt.legend() plt.subplot(132) sb.violinplot( x = 'status', y = 'axil\_nodes', data=df, marker='o', alp ha=0.5, hue = 'status') plt.xlabel('status',color = 'b',size = '13') plt.ylabel('Axil nodes',color = 'b',size = '13') plt.title("Status vs Axial\_nodes", color = 'b', size = 13) plt.legend() plt.subplot(133) sb.violinplot( x = 'status', y = 'operation\_year', data=df, marker='o', alpha=0.4, hue = 'status') plt.xlabel('status',color = 'b',size = '13') plt.ylabel('Operation year', color = 'b', size = '13') plt.title("Status vs Operation\_year",color = 'b',size = 13) plt.legend() plt.subplots\_adjust(wspace=0.5) Status vs age Status vs Operation year Status vs Axial\_nodes 90 70.0 50 80 67.5 40 70 65.0 30 60 Operation age 62.5 20 50 60.0 10 40 57.5 0 30 55.0 20 -101 1 status status status Summary of Observations 306 is the row count which indictes the number of patients and 4 is the column count which indicates the variable count 225 and 81 shows the count of patients who have survived more than 5 years and less than 5 years after the operation respectively. From the scatter plots we can infer that most the patients dont have axil\_nodes (i.e axil\_nodes =0 )and very few patients have axil\_nodes above 30 More number of patients have axil\_nodes < 5 and very less number of patients</li> having more axil\_nodes and in the plot the number of axial nodes increase the density of patients decreases gradually and patients having axial\_nodes above 30 are very few we cannot make any observation on the operation\_year vs status from the pair plot since most of the datapoints were overlapped Opeartions were happened between 1958 to 1969 The patients are aged between 30 and 83 The mean of the status is 1.26 which little bit closer to 1 than 2 .So we can say that more than half of the people lived for more than 5 years after the operation There are patients with 0 axil\_nodes The max age is 77 for the patients who lived more than 5 years Highest number of axil\_nodes found is 46 for the patients who lived more than 5 years count of patients is 225 and mean of axil\_nodes is 2.79 for the patients who lived more than 5 years The mean of axil\_nodes is 7.456 which is almost 2.5 times more than the axil\_nodes we found in the patients who lived more than 5 years so we can say that the patient having more than 46 axil\_nodes died with in 5 years after the operation The patient having age more than 77 years have not survived more than 5 years after the operation From the PDFs age vs status graph and operation\_year vs status graph didnt provide much insight because the data is overlapping in both plots.so we cant make a proper analysis By observing the values we can say that the percentage of patients with '0' axil\_nodes is 77 and with 52 axil\_nodes is 0.3 From the CDF we can observe that the patients who are having axil\_nodes greater than 30 are very few. U can clearly observe that the graph is closer to the axis as the count of axil\_nodes are increasing beyond 30 But from the axil\_nodes vs status graph we can say that the probability of patients having '0' axil\_nodes is high and the patients having axil\_nodes grater than 30 are very less From the stripp plot we can say that more patients have lived for more than 5 years 50% of patients have age less than 53 years 75% of patients have age less than 62 years I have upload the same code to my github repository Here is the link: <a href="https://github.com/bharathpreetham/EDA-on-haberman-">https://github.com/bharathpreetham/EDA-on-haberman-</a> dataset/blob/master/haberman.ipynb In [ ]: In [ ]: In [ ]: In [ ]: