Name: Akash Kulkarni Class: TE-10 Roll: 33241 Lab: DSDBAL **Problem Statement:** Perform the following operations using Python on the Air quality and Heart Diseases data sets a. Data cleaning b. Data integration c. Data transformation d. Error correcting e. Data model building import pandas as pd import numpy as np import seaborn as sns import matplotlib.pyplot as plt dataset = pd.read csv('Heart.csv') dataset.head(10) Unnamed: 0 Age Sex ChestPain RestBP Chol Fbs RestECG MaxHR ExAng Oldpeak Slope Thal AHD 0 2.3 3 0.0 1 63 1 typical 145 233 2 150 0 fixed No 1 2 67 1 asymptomatic 160 286 0 108 1.5 2 3.0 normal Yes 2 0 2 1 2.6 2 2.0 reversable 3 67 1 asymptomatic 120 229 129 Yes 3 187 3.5 37 nonanginal 130 250 3 0.0 No normal 4 5 41 0 0 2 172 1 0.0 nontypical 130 204 1.4 normal No 5 8.0 6 56 nontypical 120 236 0 178 1 0.0 normal No 6 7 62 0 asymptomatic 140 268 0 2 160 0 3.6 3 2.0 normal Yes 7 57 0 asymptomatic 120 354 163 0.6 1 0.0 normal No 8 9 63 0 2 0 1 asymptomatic 130 254 147 1.4 2 1.0 reversable Yes 10 53 140 203 155 3.1 1 asymptomatic 3 0.0 reversable Yes In [4]: dataset.isna().sum() Out[4]: Unnamed: 0 Age Sex 0 ChestPain 0 RestBP 0 Chol Fbs 0 RestECG 0 MaxHR 0 ExAng Oldpeak 0 Slope 0 4 Ca Thal AHD 0 dtype: int64 **Data Cleaning** dataset = dataset.dropna(axis=0) dataset.isnull().sum() Out[6]: Unnamed: 0 0 Sex 0 ChestPain 0 RestBP Chol 0 Fbs 0 RestECG MaxHR ExAng 0 Oldpeak 0 Slope 0 Ca Thal 0 AHD 0 dtype: int64 dataset.describe() **Unnamed:** Age Sex RestBP Chol Fbs RestECG MaxHR ExAng Oldpeak Slope **count** 297.000000 297.000000 297.000000 297.000000 297.000000 297.000000 297.000000 297.000000 297.000000 297.000000 297.000000 mean 150.673401 54.542088 0.676768 131.693603 247.350168 0.144781 0.996633 149.599327 0.326599 1.055556 1.602694 std 87.323283 9.049736 0.468500 17.762806 51.997583 0.352474 0.994914 22.941562 0.469761 1.166123 0.618187 0 1.000000 29.000000 0.000000 94.000000 126.000000 0.000000 0.000000 71.000000 0.000000 0.000000 1.000000 0 min 25% 75.000000 48.000000 0.000000 120.000000 211.000000 0.000000 0.000000 133.000000 0.000000 0.000000 1.000000 0 150.000000 56.000000 1.000000 130.000000 243.000000 0.000000 1.000000 153.000000 0.000000 0.800000 2.000000 0 **75%** 226.000000 61.000000 1.000000 140.000000 276.000000 0.000000 2.000000 166.000000 1.000000 1.600000 2.000000 1 302.000000 77.000000 1.000000 200.000000 564.000000 1.000000 2.000000 202.000000 1.000000 6.200000 3.000000 dataset.count() Out[8]: Unnamed: 0 297 Age 297 Sex 297 ChestPain 297 RestBP Chol 297 Fbs 297 RestECG 297 MaxHR 297 ExAng 297 Oldpeak 297 297 Slope Ca Thal 297 AHD 297 dtype: int64 **Data Visualizations** plt.figure(figsize=(6,4)) dataset.plot(kind='box', figsize=(20,20)) Out[31]: <AxesSubplot:> <Figure size 432x288 with 0 Axes> 0 0 100 Unnamed: 0 RestBP Chol Fbs RestECG MaxHR ExÁng Oldpeak Slope plt.figure(figsize=(10,8)) sns.distplot(dataset['Age']) plt.show() C:\Users\akash\anaconda3\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: `distplot` is a deprec ated 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). warnings.warn(msg, FutureWarning) 0.05 0.04 0.03 Density 0.02 0.01 0.00 30 40 60 70 50 plt.figure(figsize=(10,8)) sns.lmplot(x='Age', y='RestBP', data=dataset) plt.show() <Figure size 720x576 with 0 Axes> 200 180 160 140 120 100 30 40 50 60 70 Age plt.figure(figsize=(10,8)) sns.lmplot(x='Age', y='Chol', data=dataset) plt.show() <Figure size 720x576 with 0 Axes> 500 400 g 300 60 70 plt.figure(figsize=(10,8)) sns.heatmap(dataset.corr(), annot=True) - 1.0 Unnamed: 0 - 1 0.0093 -0.088 -0.022 -0.084 -0.052 -0.14 -0.12 -0.0027 -0.11 -0.032 0.049 Age -0.0093 - 0.8 -0.088 -0.092 -0.066 0.039 0.034 0.14 0.033 0.11 -0.066 - 0.6 RestBP - -0.022 0.29 0.18 -0.084 -0.2 0.13 0.013 0.17 -7.5e-05 0.059 0.039 -0.0092 0.12 Chol -- 0.4 0.039 0.18 0.013 0.069 -0.00780.000890.0083 0.048 Fbs - -0.052 0.13 0.034 0.15 0.17 0.069 -0.072 0.082 0.11 0.15 RestECG -- 0.2 -0.06 -0.049-7.5e-050.0078-0.072 -0.35 -0.39 MaxHR -0.38 0.14 0.067 0.059-0.000890.082 -0.38 ExAng -0.0027 0.096 0.25 0.15 - 0.0 0.039 0.0083 0.11 0.29 -0.35 0.29 Oldpeak · -0.11- -0.2 0.033 0.12 -0.0092 0.048 -0.032 0.16 -0.39 0.25 1 0.11 -0.27 0.049 0.092 0.098 0.15 0.15 0.29 Slope Ö In [14]: sns.pairplot(dataset, kind='reg', plot kws={'line kws': {'color':'cyan'}}) Out[14]: <seaborn.axisgrid.PairGrid at 0x1edbafe44c0> **Data Transformation** print(dataset.ChestPain.unique()) print(dataset.Thal.unique()) print(dataset.AHD.unique()) ['typical' 'asymptomatic' 'nonanginal' 'nontypical'] ['fixed' 'normal' 'reversable'] ['No' 'Yes'] heart\_encoding = pd.get\_dummies(dataset[['ChestPain', 'Thal', 'AHD']]) heart\_final = pd.concat([dataset, heart\_encoding],1) heart\_final = heart\_final.drop(['ChestPain', 'Thal', 'AHD'], axis = 1) heart\_final.head(10) **Unnamed:** Sex RestBP Chol Fbs RestECG MaxHR ExAng Oldpeak ... Ca ChestPain\_asymptomatic ChestPain\_nonanginal Chest Age 0 150 0 0 0 1 63 1 145 233 1 2 2.3 ... 0.0 1 2 67 160 286 108 1.5 ... 3.0 2 3 67 1 120 229 0 2 129 1 2.6 ... 2.0 1 0 ... 0.0 4 37 1 130 250 0 187 3.5 4 0 0 0 5 41 0 130 204 2 172 0 1.4 ... 0.0 5 6 56 120 236 0 178 0.8 ... 0.0 6 7 0 0 2 0 0 62 140 268 160 3.6 ... 2.0 1 8 57 0 120 354 0 0 163 0.6 ... 0.0 8 9 0 0 0 63 1 130 254 2 147 1.4 ... 1.0 1 155 10 53 140 203 3.1 ... 0.0 10 rows × 21 columns In [18]: heart\_final.Sex.value\_counts() 201 96 Name: Sex, dtype: int64 In [19]: pd.crosstab(heart\_final.AHD\_Yes,heart\_final.Sex) Sex 0 1 AHD\_Yes **0** 71 89 **1** 25 112 heart\_final.columns Out[20]: Index(['Unnamed: 0', 'Age', 'Sex', 'RestBP', 'Chol', 'Fbs', 'RestECG', 'MaxHR', 'ExAng', 'Oldpeak', 'Slope', 'Ca', 'ChestPain\_asymptomatic', 'ChestPain\_nonanginal', 'ChestPain\_nontypical', 'ChestPain\_typical', 'Thal\_fixed', 'Thal\_normal', 'Thal\_reversable', 'AHD\_No', 'AHD\_Yes'], dtype='object') heart\_final.dtypes Out[21]: Unnamed: 0 int64 int64 Age int64 RestBP int64 Chol int64 int64 RestECG int64 MaxHR int64 ExAng int64 Oldpeak float64 Slope int64 float64 uint8 ChestPain\_asymptomatic ChestPain nonanginal uint8 ChestPain nontypical uint8 ChestPain\_typical uint8 Thal\_fixed uint8 Thal\_normal uint8 Thal reversable uint8 AHD No uint8 AHD Yes uint8 dtype: object df = heart final.drop('AHD Yes', axis=1) df norm = (df-df.min())/(df.max()-df.min())df\_norm = pd.concat((df\_norm, heart\_final.AHD\_Yes), 1) df norm.head(10) **Unnamed:** Age Sex RestBP Chol Fbs RestECG MaxHR ExAng Oldpeak ... Ca ChestPain\_asymptomatic ChestPain\_no 1.0 0.603053 0 0.000000 0.708333 1.0 0.481132 0.244292 1.0 0.0 0.370968 ... 0.000000 0.0 1 0.003322 0.791667 1.0 0.622642 0.365297 1.0 0.282443 1.0 0.241935 ... 1.000000 1.0 1.0 0.245283 0.235160 0.006645 0.791667 1.0 0.442748 1.0 0.419355 ... 0.666667 2 1.0 1.0 0.339623 0.283105 3 0.009967 0.166667 0.0 0.885496 0.0 0.564516 ... 0.000000 0.0 4 0.013289 0.250000 0.0 0.339623 0.178082 1.0 0.770992 0.0 0.225806 ... 0.000000 0.0 0.016611 0.562500 5 1.0 0.245283 0.251142 0.0 0.0 0.816794 0.0 0.129032 ... 0.000000 0.0 0.019934 0.687500 0.0 0.433962 0.324201 1.0 0.679389 0.0 0.580645 ... 0.666667 6 1.0 0.0 0.245283 0.520548 7 0.023256 0.583333 0.0 0.702290 1.0 0.096774 ... 0.000000 1.0 0.026578 0.708333 1.0 0.339623 0.292237 0.0 1.0 0.580153 0.0 0.225806 ... 0.333333 1.0 0.029900 0.500000 1.0 0.433962 0.175799 1.0 1.0 0.641221 1.0 0.500000 ... 0.000000 10 rows × 21 columns In [24]: X = df\_norm.drop(['AHD\_Yes', 'Unnamed: 0'], axis=1) Y = df norm.AHD Yes**Train Test Split** from sklearn.model\_selection import train\_test\_split X\_train,X\_test,y\_train,y\_test=train\_test\_split(X,Y,test\_size=0.2 ,random\_state=1) **KNN Modelling** from sklearn.neighbors import KNeighborsClassifier classifier = KNeighborsClassifier(n neighbors=5, metric="minkowski", p=2) classifier.fit(X\_train, y\_train) Out[26]: KNeighborsClassifier() y\_pred = classifier.predict(X\_test) y\_pred 0, 0, 0, 0, 1, 1, 0, 0, 0, 1, 0, 1, 1, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 1, 0, 1, 0, 1, 0, 1, 0, 1, 1, 1, 1, 1, 1, 0, 1], dtype=uint8) from sklearn.metrics import confusion\_matrix conf mat = confusion matrix(y test, y pred) sns.heatmap(conf\_mat, square=True, annot=True, cmap='Blues', fmt='d', cbar=False) Out[28]: <AxesSubplot:> 30 from sklearn.metrics import accuracy score accuracy\_score(y\_test,y\_pred) Out[30]: 0.983333333333333333 Conclusion: Accuracy of trained model is 98.34%