

Jeevan Kumar Banoth - 02105145

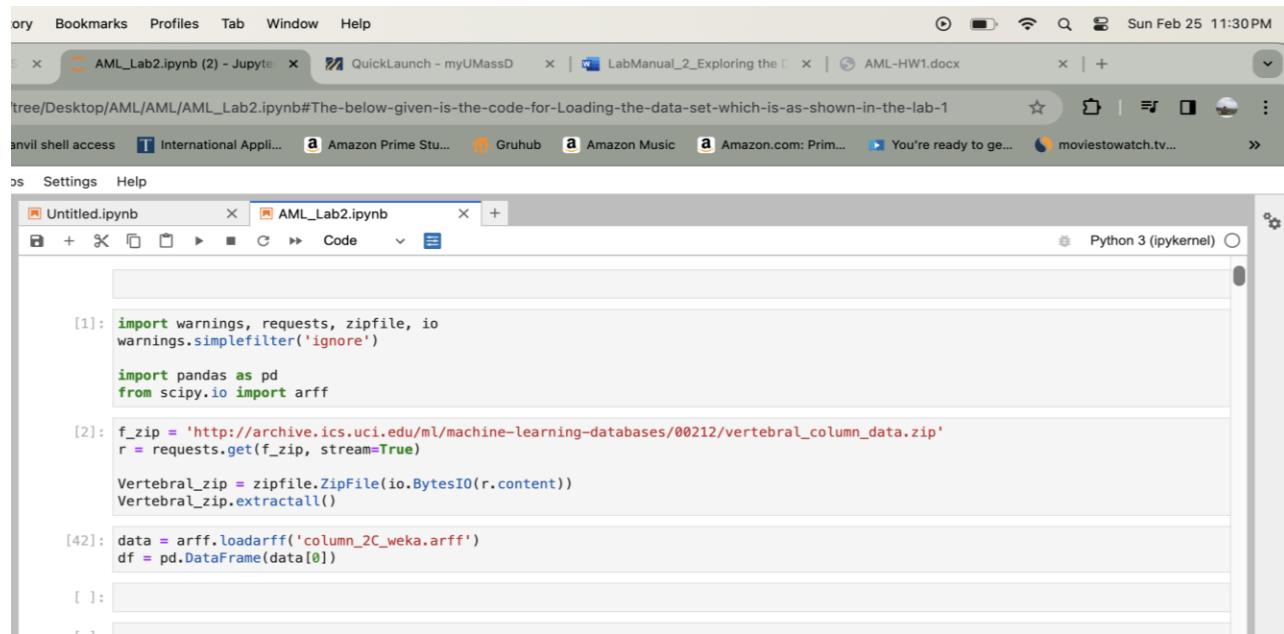
Advanced machine learning

Homework –2:

Lab setup:

The below given is the code for Loading the data set, which is as shown in the lab 1. Because the solution is split across several labs in this module, we have imported that data set I have imported that as instructed in the Lab manual

Importing the data:



A screenshot of a web browser window showing a Jupyter Notebook interface. The browser tabs include 'AML_Lab2.ipynb (2) - Jupyter', 'QuickLaunch - myUMassD', 'LabManual_2_Exploring the...', 'AML-HW1.docx', and several other unrelated tabs like 'Amazon Prime Stu...', 'Gruhub', 'Amazon Music', etc. The Jupyter notebook cell contains the following Python code:

```
[1]: import warnings, requests, zipfile, io
warnings.simplefilter('ignore')

import pandas as pd
from scipy.io import arff

[2]: f_zip = 'http://archive.ics.uci.edu/ml/machine-learning-databases/00212/vertebral_column_data.zip'
r = requests.get(f_zip, stream=True)

Vertebral_zip = zipfile.ZipFile(io.BytesIO(r.content))
Vertebral_zip.extractall()

[42]: data = arff.loadarff('column_2C_weka.arff')
df = pd.DataFrame(data[0])

[ ]:
```

So now the data has been imported we can start Lab-2 from the steps below.

Step 1: Exploring the data

First, we are using “shape” to determine the number of rows and columns in the dataset. As we have given “df” to the name of dataset, we are using “df.shape” command for determining the number of rows and columns in the dataset.

The screenshot shows a Jupyter Notebook interface within a Chrome browser. The notebook has two tabs open: 'Untitled.ipynb' and 'AML_Lab2.ipynb'. The 'Untitled.ipynb' tab is active, displaying the following code and output:

```
[41]: df.shape
[41]: (310, 7)

[48]: print(df.dtypes)
pelvic_incidence    float64
pelvic_tilt         float64
lumbar_lordosis_angle  float64
sacral_slope        float64
pelvic_radius       float64
degree_spondylolisthesis float64
class               object
dtype: object
```

As we can see here it shows the number of rows as 310 and the number of columns 7.

Here, dtypes shows the types of data that we have in the dataset, Here we are using the df.dtypes command, which shows types of data present in that.

The screenshot shows a Jupyter Notebook interface within a Chrome browser. The notebook has two tabs open: 'Untitled.ipynb' and 'AML_Lab2.ipynb'. The 'Untitled.ipynb' tab is active, displaying the following code and output:

```
[47]: print(df.columns)
Index(['pelvic_incidence', 'pelvic_tilt', 'lumbar_lordosis_angle',
       'sacral_slope', 'pelvic_radius', 'degree_spondylolisthesis', 'class'],
      dtype='object')

[46]: df
[46]:
```

Below the code, the data is displayed as a table:

	pelvic_incidence	pelvic_tilt	lumbar_lordosis_angle	sacral_slope	pelvic_radius	degree_spondylolisthesis	class
0	63.027817	22.552586	39.609117	40.475232	98.672917	-0.254400	b'Abnormal'
1	39.056951	10.060091	25.015378	28.995960	114.405425	4.564259	b'Abnormal'
2	68.832021	22.218482	50.092194	46.613539	105.985135	-3.530317	b'Abnormal'
3	69.297008	24.652878	44.311238	44.644130	101.868495	11.211523	b'Abnormal'
4	49.712859	9.652075	28.317406	40.060784	108.168725	7.918501	b'Abnormal'
...
305	47.903565	13.616688	36.000000	34.286877	117.449062	-4.245395	b'Normal'
306	53.936748	20.721496	29.220534	33.215251	114.365845	-0.421010	b'Normal'
307	61.446597	22.694968	46.170347	38.751628	125.670725	-2.707880	b'Normal'
308	45.252792	8.693157	41.583126	36.559635	118.545842	0.214750	b'Normal'
309	33.841641	5.073991	36.641233	28.767649	123.945244	-0.199249	b'Normal'

310 rows × 7 columns

- "df" command is used here, which shows all data dataset rows and columns

- Columns command is used for showing the number of columns present in the dataset.
- So, here we are using “df.columns” command which will represent the number of columns.

→ “df” command is used here, which shows all data dataset rows and columns.

The screenshot shows a Jupyter Notebook interface within a Chrome browser. The notebook has two cells displayed:

```
[50]: df.describe()
[50]:
```

	pelvic_incidence	pelvic_tilt	lumbar_lordosis_angle	sacral_slope	pelvic_radius	degree_spondylolisthesis
count	310.000000	310.000000	310.000000	310.000000	310.000000	310.000000
mean	60.496653	17.542822	51.930930	42.953831	117.920655	26.296694
std	17.236520	10.008330	18.554064	13.423102	13.317377	37.559027
min	26.147921	-6.554948	14.000000	13.366931	70.082575	-11.058179
25%	46.430294	10.667069	37.000000	33.347122	110.709196	1.603727
50%	58.691038	16.357689	49.562398	42.404912	118.268178	11.767934
75%	72.877696	22.120395	63.000000	52.695888	125.467674	41.287352
max	129.834041	49.431864	125.742385	121.429566	163.071041	418.543082

```
[12]: df['pelvic_incidence'].describe()
[12]:
```

	pelvic_incidence	dtype
count	310.000000	
mean	60.496653	
std	17.236520	
min	26.147921	
25%	46.430294	
50%	58.691038	
75%	72.877696	
max	129.834041	

→ “df.describe()” command used to describe their statistics of entire data set.

→ For looking at their statistics of each row we can use the name of that row and the describe () function so that it will display the actual functionalities of that respective row.

→ Below are the few different codes which describe their statistics of respective rows.

```

[55]: df['sacral_slope'].describe()

[55]: count    310.000000
      mean     42.953831
      std      13.423102
      min     13.366931
      25%    33.347122
      50%    42.404912
      75%    52.695888
      max    121.429566
      Name: sacral_slope, dtype: float64

[56]: df.describe()

[56]:   pelvic_incidence  pelvic_tilt  lumbar_lordosis_angle  sacral_slope  pelvic_radius  degree_spondylolisthesis
      count  310.000000  310.000000  310.000000  310.000000  310.000000
      mean   60.496653  17.542822  51.930930  42.953831  117.920655  26.296694
      std    17.236520  10.008330  18.554064  13.423102  13.317377  37.559027
      min    26.147921  -6.554948  14.000000  13.366931  70.082575  -11.058179
      25%   46.430294  10.667069  37.000000  33.347122  110.709196  1.603727
      50%   58.691038  16.357689  49.562398  42.404912  118.268178  11.767934
      75%   72.877696  22.120395  63.000000  52.695888  125.467674  41.287352
      max   129.834041  49.431864  125.742385  121.429566  163.071041  418.543082

```

→ By changing the code as follows we can get their statistics.

Features with potential issues:

`degree_spondylolisthesis`:

Large standard deviation (37.5) compared to the mean (2.6) suggests a wide spread of values.

Minimum value (-11.06) and maximum value (418.54) indicate potential outliers.

Features with possible correlations:

`pelvic_incidence`, `pelvic_tilt`, `sacral_slope`: These features are all related to the pelvic alignment and might show correlations.

`pelvic_tilt`, `lumbar_lordosis_angle`: Pelvic tilt and lumbar lordosis often compensate for each other, potentially showing negative correlation.

Distribution visualizations:

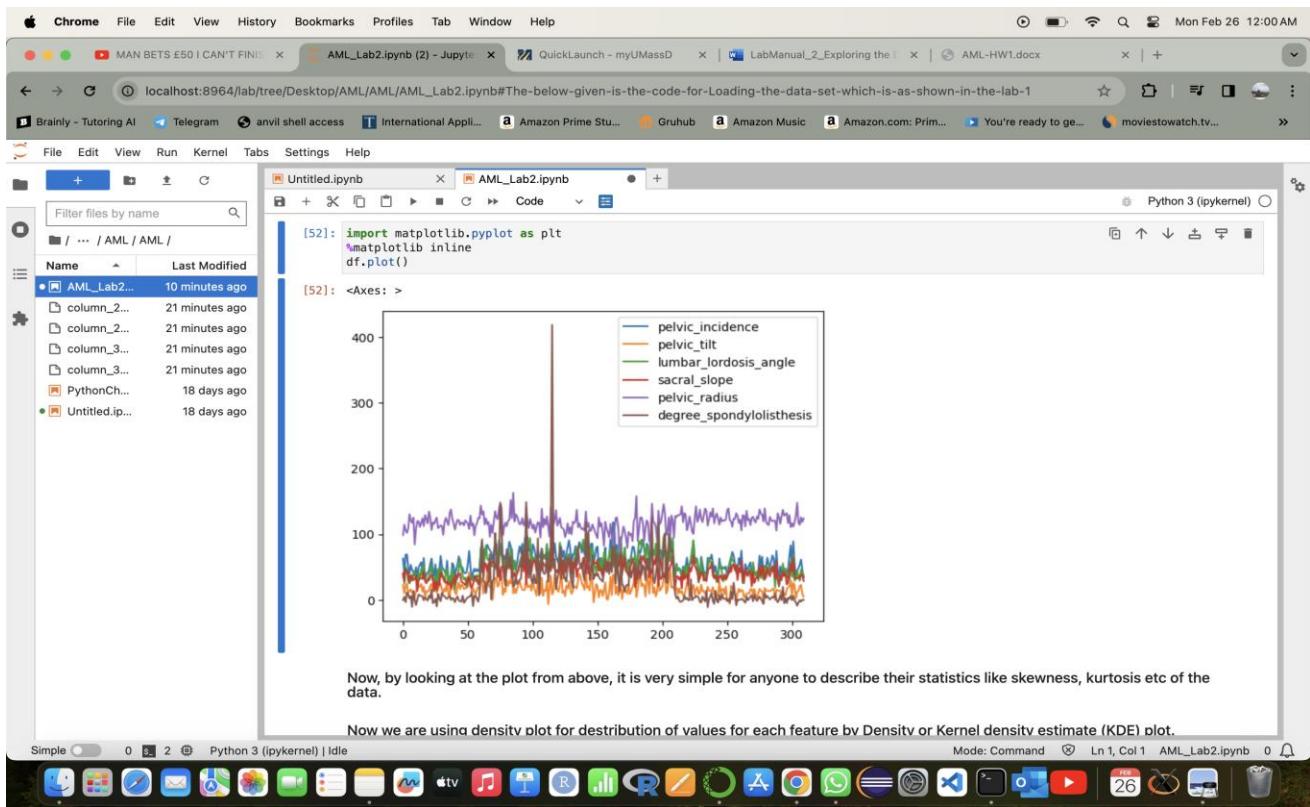
Histograms or boxplots can better visualize distributions and identify outliers.

Correlation matrix:

Calculating and visualizing the correlation matrix would reveal potential relationships between features.

As it is not always easy to make the observations based on their numbers, so here we are going with plotting them based on the given data.

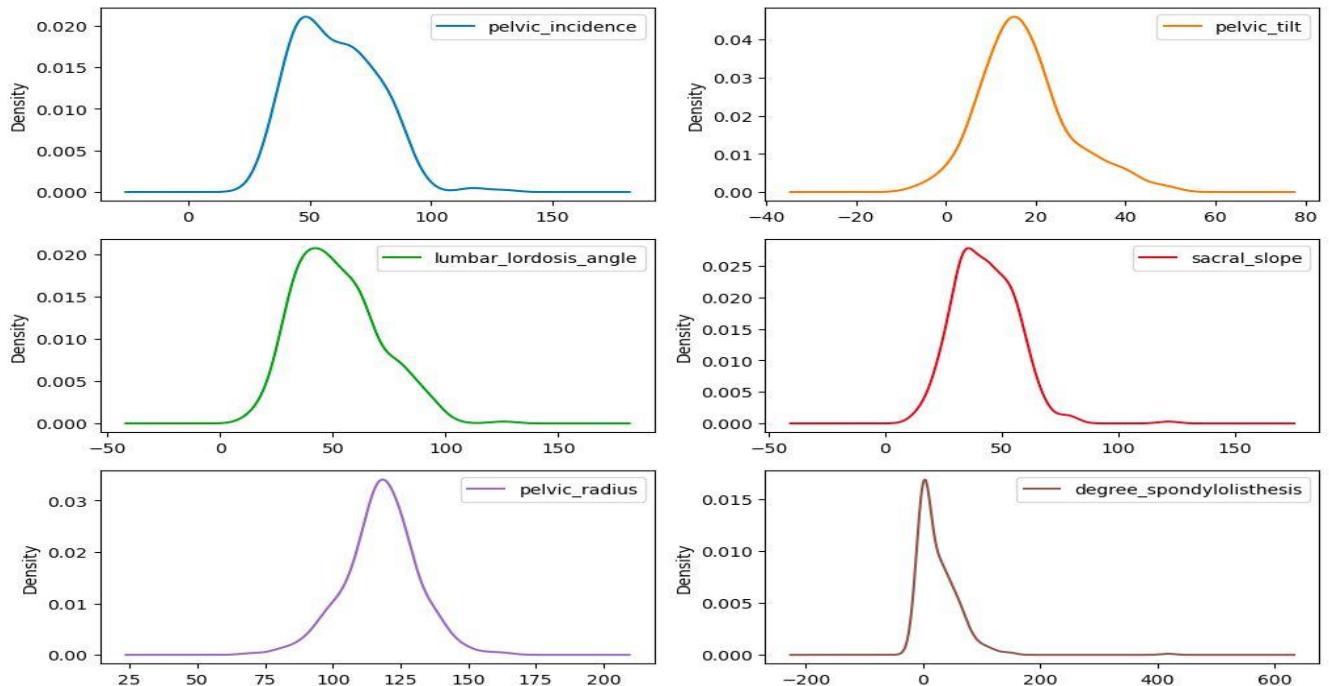
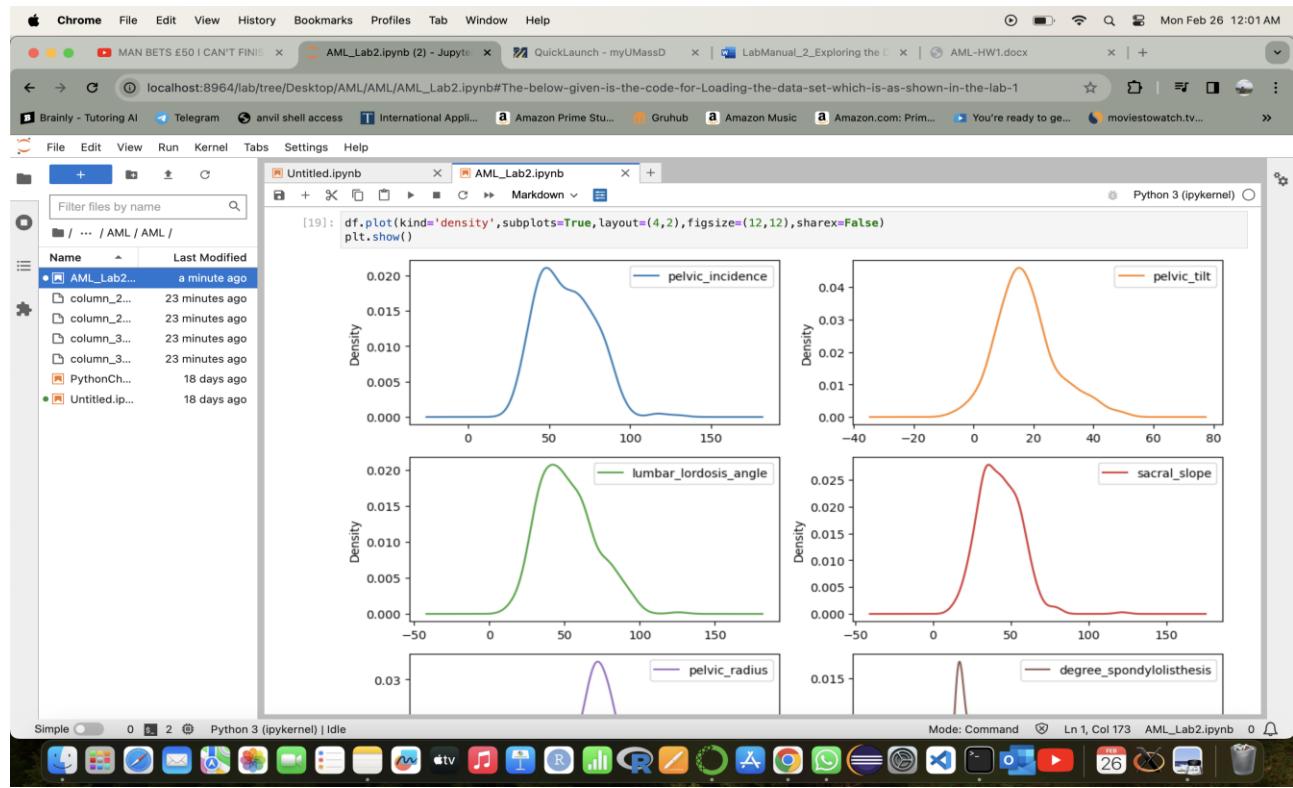
The below code is for plotting the dataset:



Now, by looking at the plot from above, it is very simple for anyone to describe their statistics like skewness, kurtosis etc of the data.

Now we are using density plot for distribution of values for each feature by Density or Kernel density estimate (KDE) plot.

By performing this code, we can now get the different plots for each function or row as below. This makes it easier for one to describe or differentiate them.



Here are a few observations:

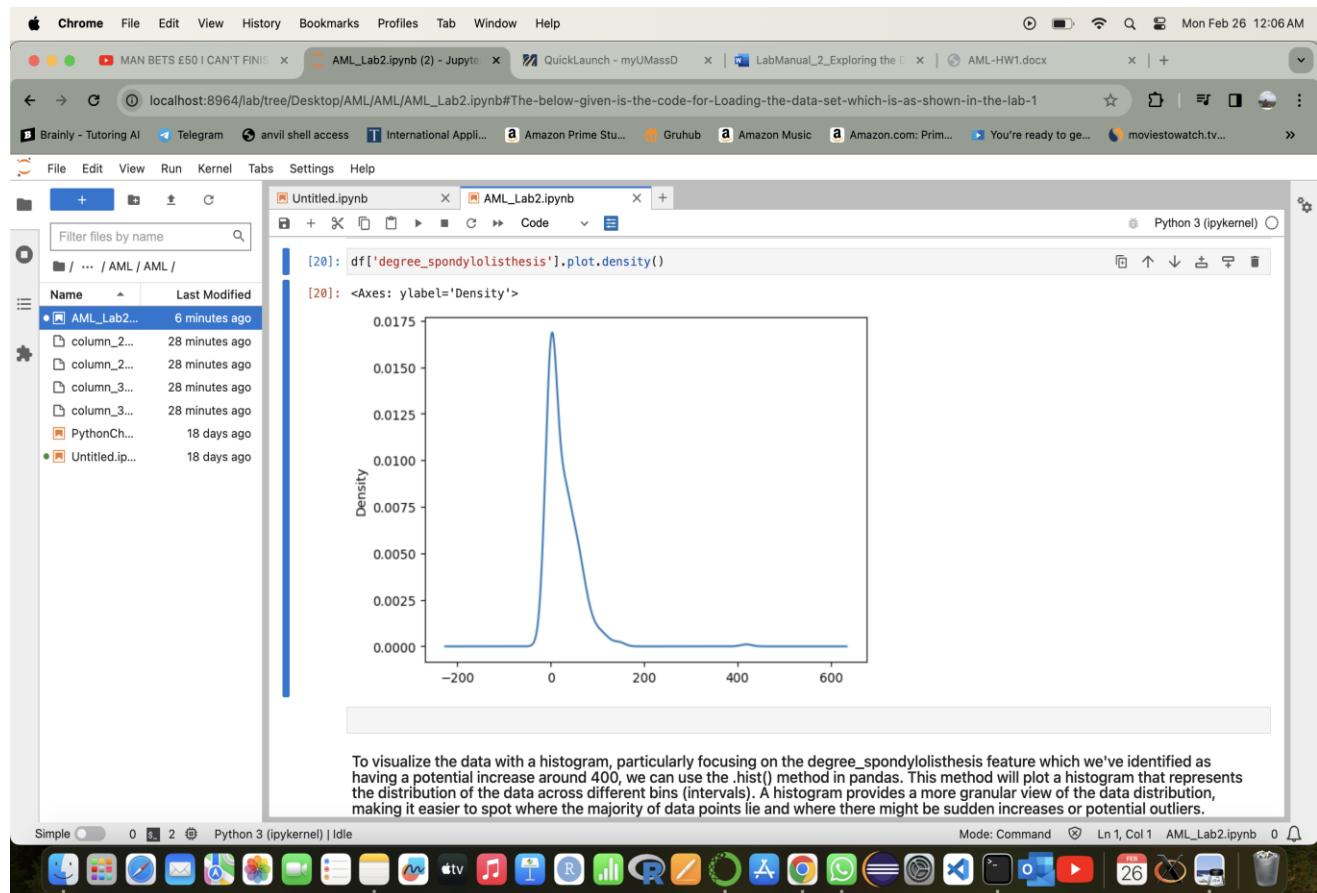
The distribution of pelvic tilt appears to be bimodal, with two peaks at around -10 and 20 degrees. This suggests that there may be two distinct groups of people in the dataset, with one group having a more anteriorly tilted pelvis and the other group having a more posteriorly tilted pelvis.

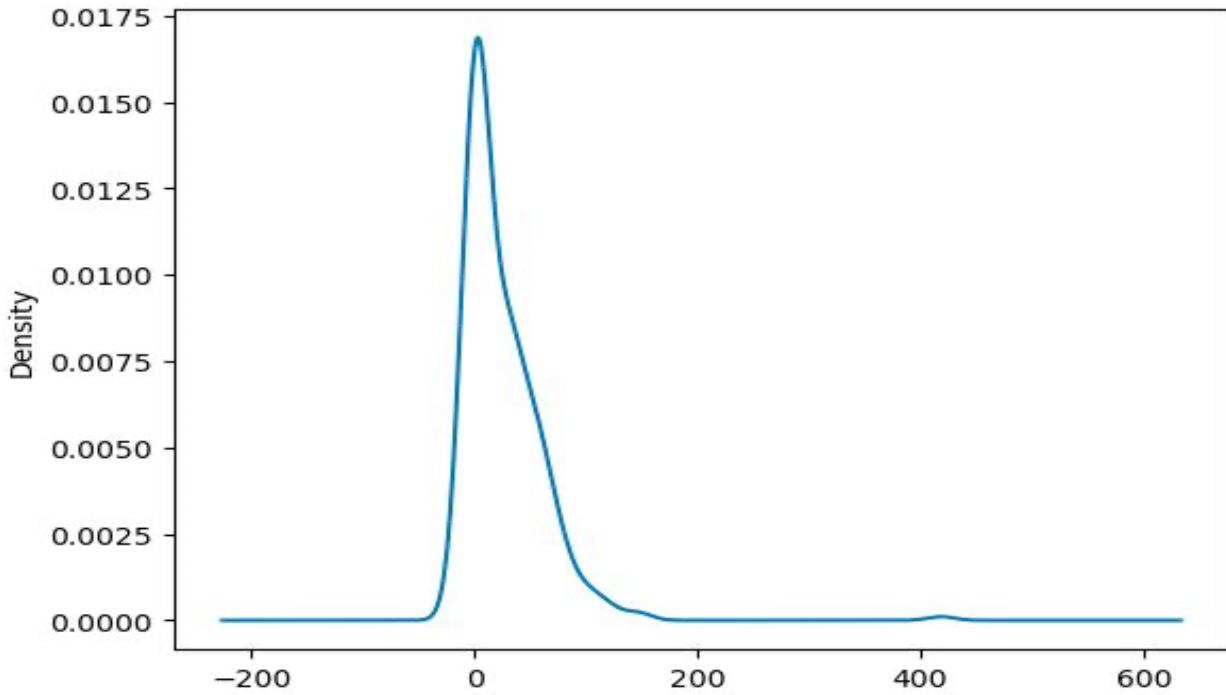
The distribution of lumbar lordosis angle appears to be skewed to the right, with a tail extending towards higher values. This suggests that most people have a moderate amount of lordosis, but some people have a very large amount of lordosis.

The distribution of degree of spondylolisthesis appears to be centered around zero, with a few outliers on either side. This suggests that most people do not have spondylolisthesis, but a small number of people do have it.

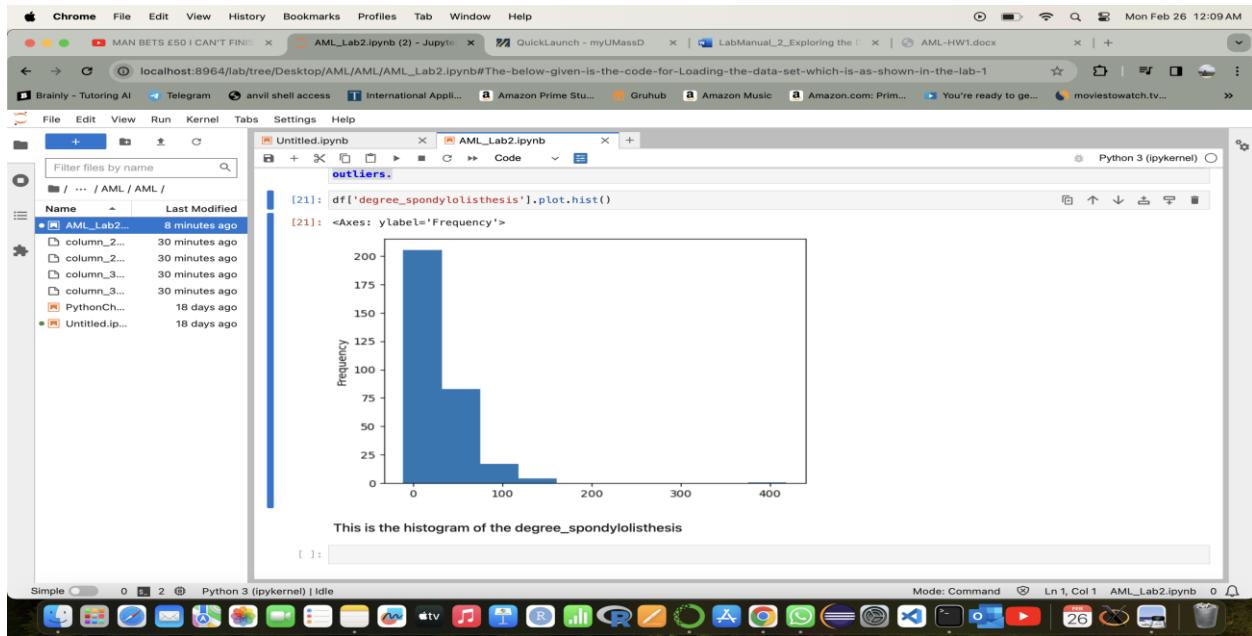
Investigating degree spondylolisthesis:

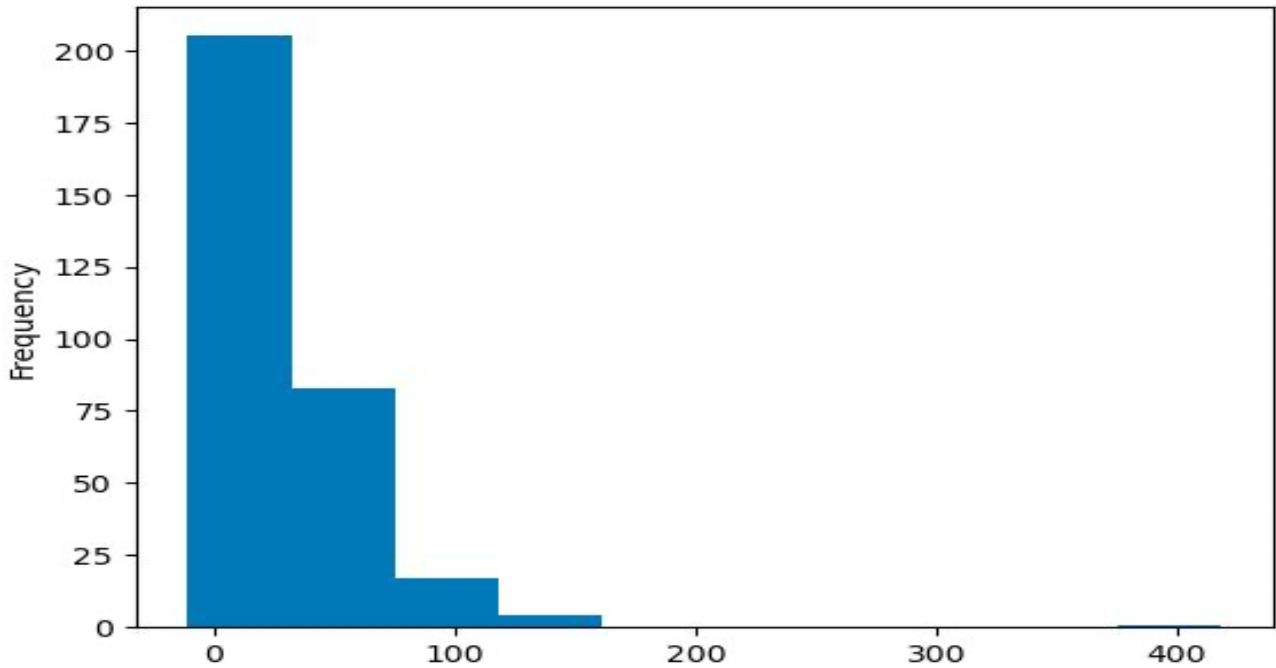
Now we are investigating degree spondylolisthesis, we are starting with density plot which shows the distribution of values





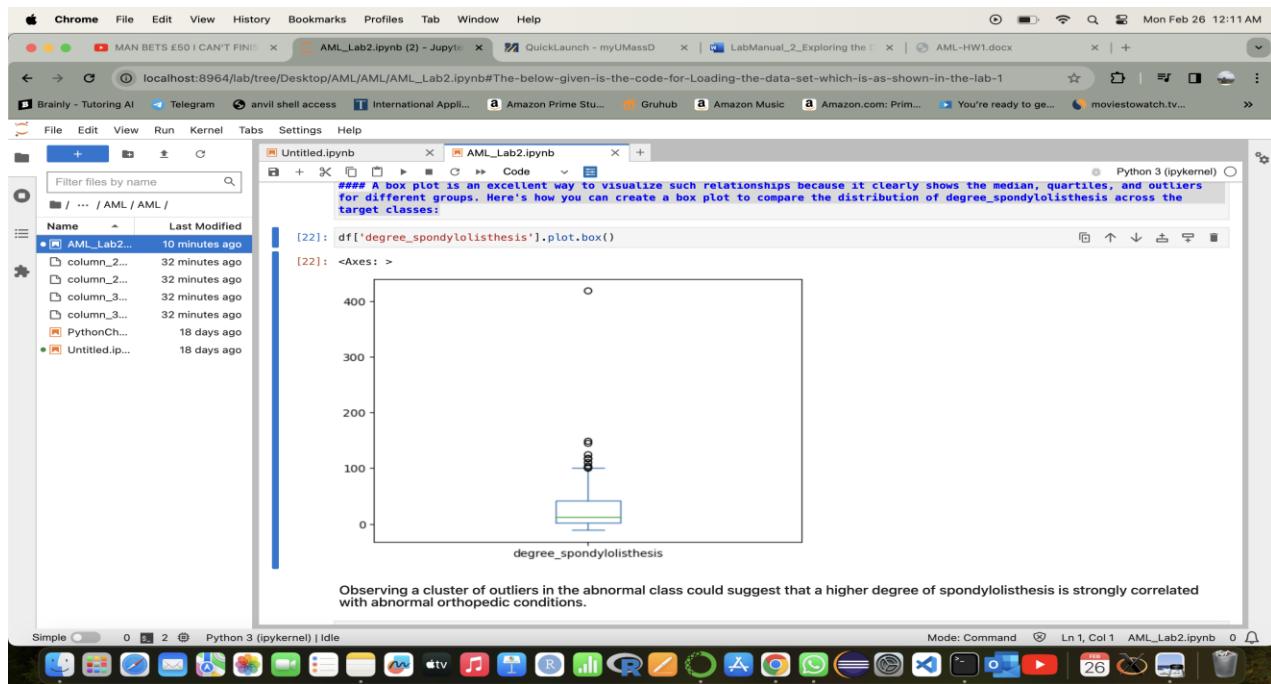
To visualize the data with a histogram, particularly focusing on the degree spondylolisthesis feature which we've identified as having a potential increase around 400, we can use the `.hist()` method in pandas. This method will plot a histogram that represents the distribution of the data across different bins (intervals). A histogram provides a more granular view of the data distribution, making it easier to spot where most data points lie and where there might be sudden increases or potential outliers.





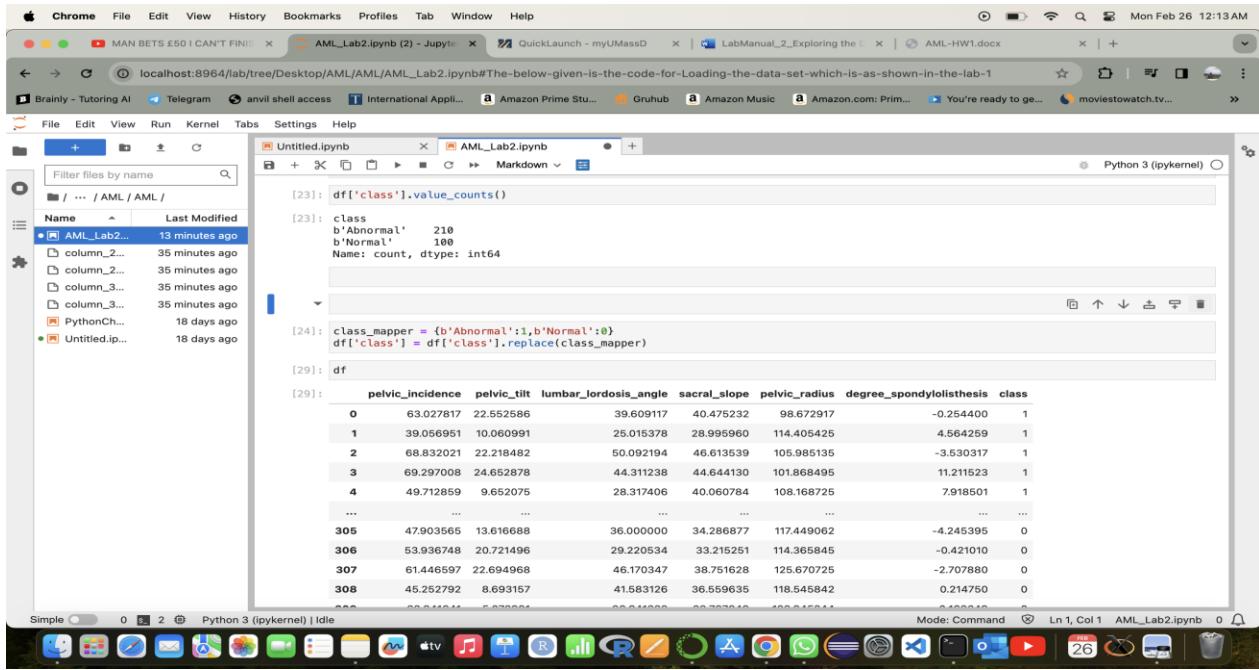
→ This is the histogram of the degree_spondylolisthesis

A box plot is an excellent way to visualize such relationships because it clearly shows the median, quartiles, and outliers for different groups. Here's how you can create a box plot to compare the distribution of degree_spondylolisthesis across the target classes:



Observing a cluster of outliers in the abnormal class could suggest that a higher degree of spondylolisthesis is strongly correlated with abnormal orthopedic conditions.

Analyzing the target:



The screenshot shows a Jupyter Notebook running in a Chrome browser window. The notebook has two tabs open: 'Untitled.ipynb' and 'AML_Lab2.ipynb'. The 'Untitled.ipynb' tab is active, displaying the following code and output:

```
[23]: df['class'].value_counts()  
[23]: class  
b'Abnormal'    210  
b'Normal'     100  
Name: count, dtype: int64  
  
[24]: class_mapper = {b'Abnormal':1,b'Normal':0}  
df['class'] = df['class'].replace(class_mapper)  
  
[29]: df
```

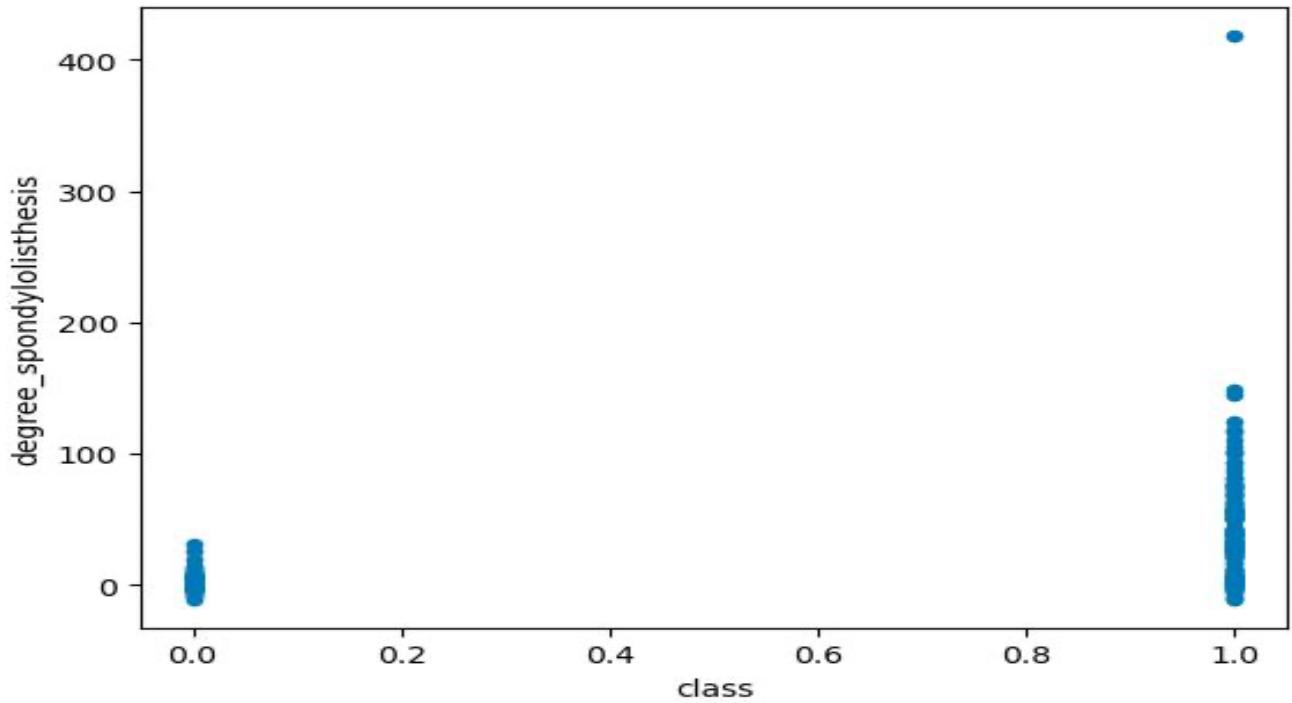
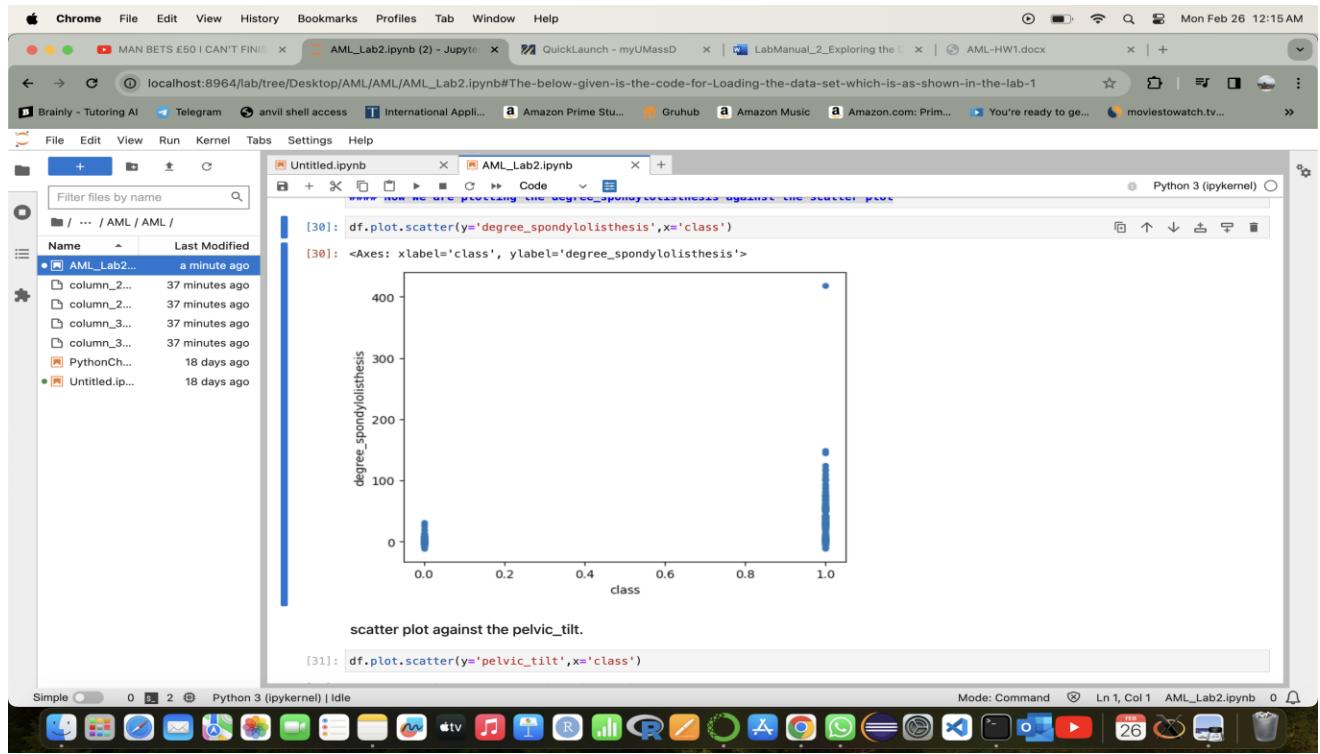
The output of the last cell shows a DataFrame with 308 rows and 8 columns. The columns are: pelvic_incidence, pelvic_tilt, lumbar_lordosis_angle, sacral_slope, pelvic_radius, degree_spondylolisthesis, and class. The 'class' column contains binary values 1 and 0, corresponding to the 'Abnormal' and 'Normal' categories respectively.

	pelvic_incidence	pelvic_tilt	lumbar_lordosis_angle	sacral_slope	pelvic_radius	degree_spondylolisthesis	class
0	63.027817	22.552586	39.609117	40.475232	98.672917	-0.254400	1
1	39.056951	10.060991	25.015378	28.995960	114.405425	4.564259	1
2	68.832021	22.218482	50.092194	46.613539	105.985135	-3.530317	1
3	69.297008	24.652878	44.311238	44.64130	101.868495	11.211523	1
4	49.712859	9.652075	28.317406	40.060784	108.168725	7.918501	1
...
305	47.903565	13.616688	36.000000	34.286877	117.449062	-4.245395	0
306	53.936748	20.721496	29.220534	33.215251	114.365845	-0.421010	0
307	61.446597	22.694968	46.170347	38.751628	125.670725	-2.707880	0
308	45.252792	8.693157	41.583126	36.559635	118.545842	0.214750	0

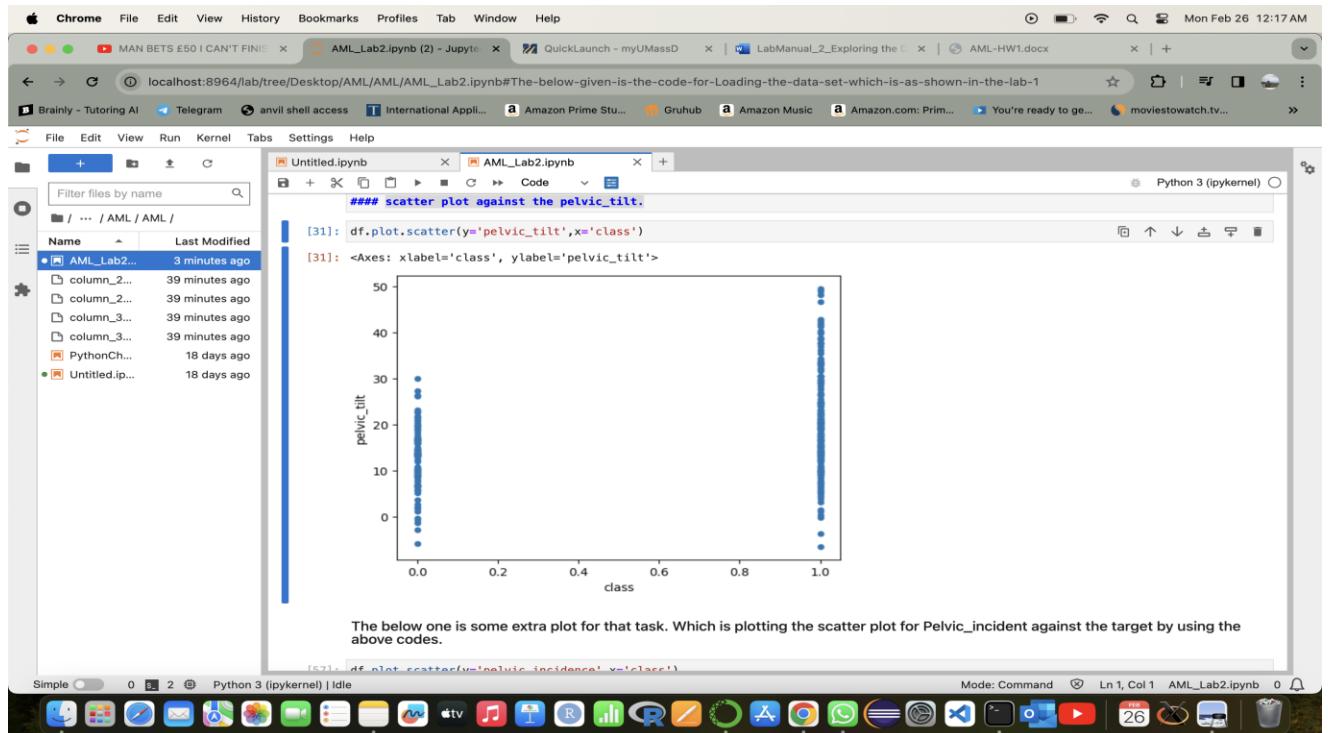
The output of `df['class'].value_counts()` shows the distribution of classes in your dataset, indicating there are 210 cases labeled as b'Abnormal' and 100 cases labeled as b'Normal'. This information is useful for understanding the balance of your dataset.

To convert the class column from binary strings to numeric values, we've defined a mapping dictionary named `class_mapper` where b'Abnormal' is mapped to 1 and b'Normal' is mapped to 0. This conversion is necessary because most machine learning algorithms work with numeric data, and having your target variable in a numeric format will facilitate model training and evaluation.

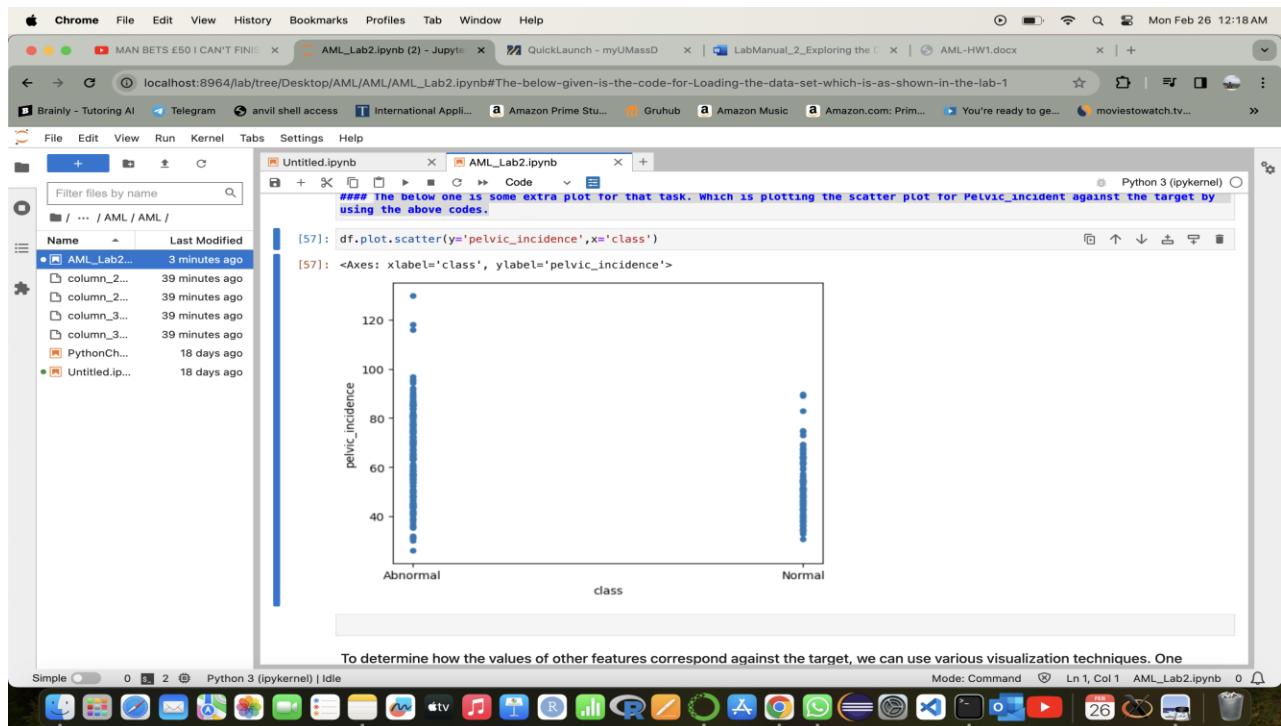
Now we are plotting the `degree_spondylolisthesis` against the scatter plot



scatter plot against the pelvic_tilt.



The below one is some extra plot for that task. Which is plotting the scatter plot for Pelvic_incident against the target by using the above codes.

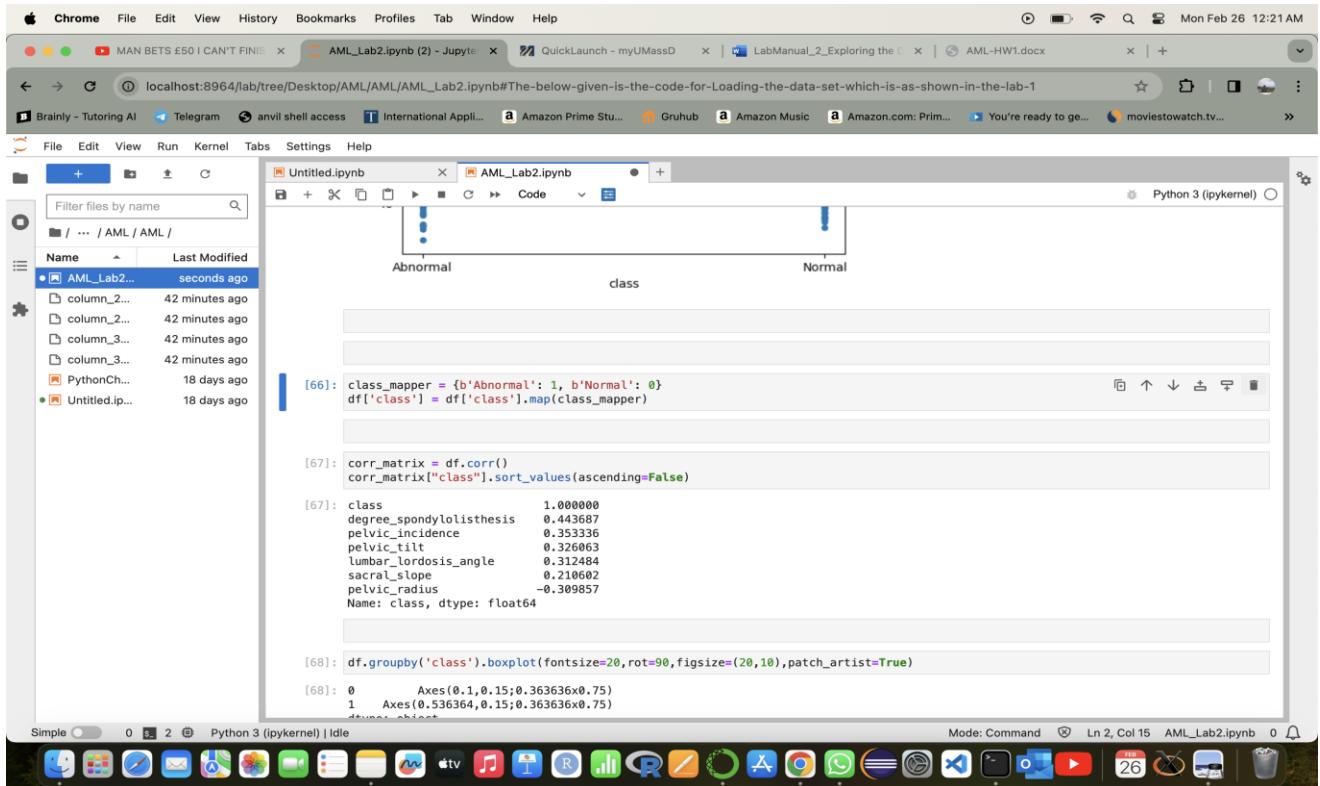


To determine how the values of other features correspond against the target, we can use various visualization techniques. One effective approach is to create scatter plots for each feature against

the target, and use color coding to differentiate between the classes. Additionally, box plots can offer insights into the distribution of feature values across different classes.

Visualizing multiple variables:

First, let's create a correlation matrix and look at how each feature correlates with the 'class'. This will give us a numerical understanding of relationships.



The screenshot shows a Jupyter Notebook interface running on a Mac OS X desktop. The browser tab bar includes 'MAN BETS £50 I CAN'T FINISH...', 'localhost:8964/lab/tree/Desktop/AML/AML_Lab2.ipynb (2) - Jupyter...', 'QuickLaunch - myUMassD...', 'LabManual_2_Exploring the...', and 'AML-HW1.docx'. The notebook sidebar shows files like 'Untitled.ipynb' and 'AML_Lab2.ipynb'. The main area displays Python code and its output. The code includes mapping 'Abnormal' to 1 and 'Normal' to 0, creating a correlation matrix, and generating a boxplot grouped by 'class'. The output shows the correlation matrix with 'class' having a value of 1.00000 and other features like 'degree_spondylolisthesis' and 'pelvic_radius' having both positive and negative correlations.

```
[66]: class_mapper = {b'Abnormal': 1, b'Normal': 0}
df['class'] = df['class'].map(class_mapper)

[67]: corr_matrix = df.corr()
corr_matrix["class"].sort_values(ascending=False)

[68]: class          1.00000
degree_spondylolisthesis    0.443687
pelvic_incidence           0.353336
pelvic_tilt                 0.326063
lumbar_lordosis_angle       0.312484
sacral_slope                  0.210602
pelvic_radius                -0.309857
Name: class, dtype: float64

[69]: df.groupby('class').boxplot(fontsize=20, rot=90, figsize=(20,10), patch_artist=True)

[70]: 0      Axes(0.1,0.15;0.363636x0.75)
1      Axes(0.536364,0.15;0.363636x0.75)
```

This code snippet prints out the correlation coefficients between each feature and the 'class', sorted in descending order. A positive value indicates a positive correlation, meaning as the feature value increases, the likelihood of the class being 'Abnormal' (1) increases. Conversely, a negative value indicates an inverse relationship.

Now we have created a correlation of matrices by using the corr function for the entire dataset.

Chrome File Edit View History Bookmarks Profiles Tab Window Help

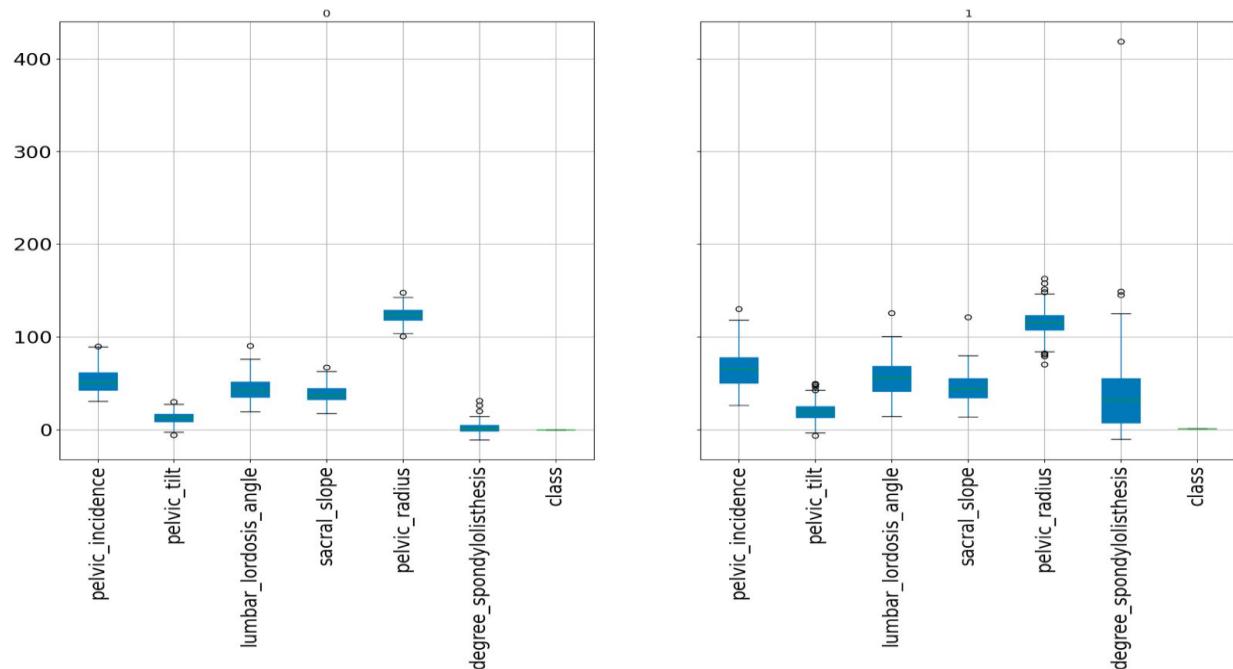
localhost:8964/lab/tree/Desktop/AML/AML_Lab2.ipynb#The-below-given-is-the-code-for>Loading-the-data-set-which-is-as-shown-in-the-lab-1

Untitled.ipynb AML_Lab2.ipynb QuickLaunch - myUMassD LabManual_2_Exploring the ... AML-HW1.docx

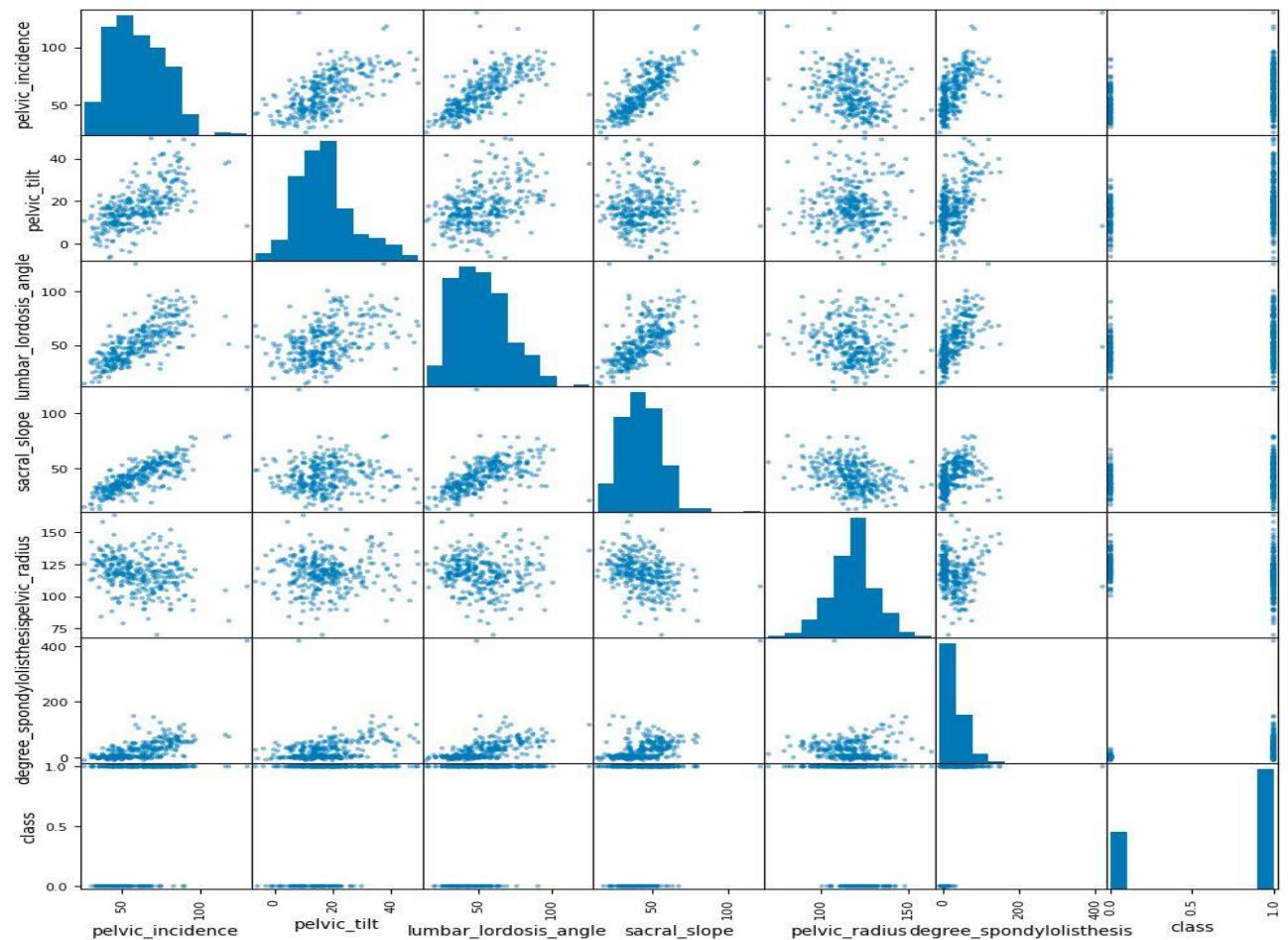
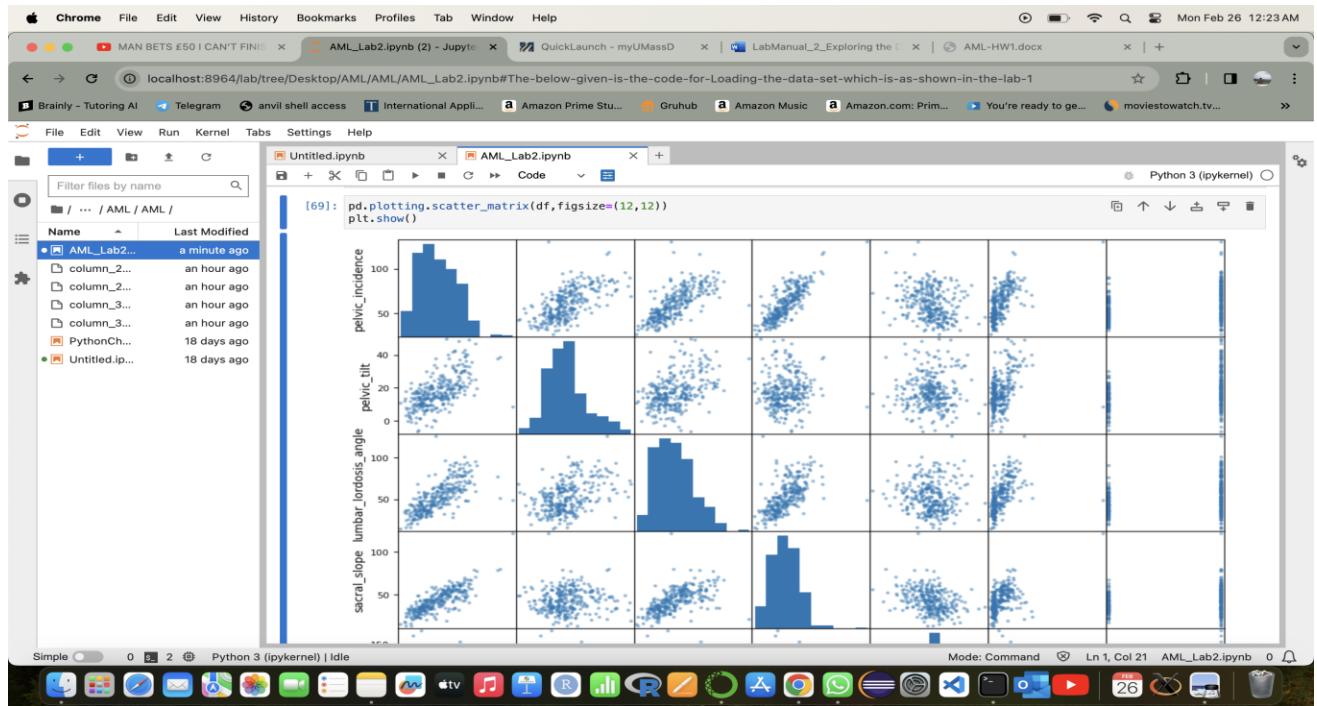
File Edit View Run Kernel Tabs Settings Help

Untitled.ipynb AML_Lab2.ipynb Python 3 (ipykernel)

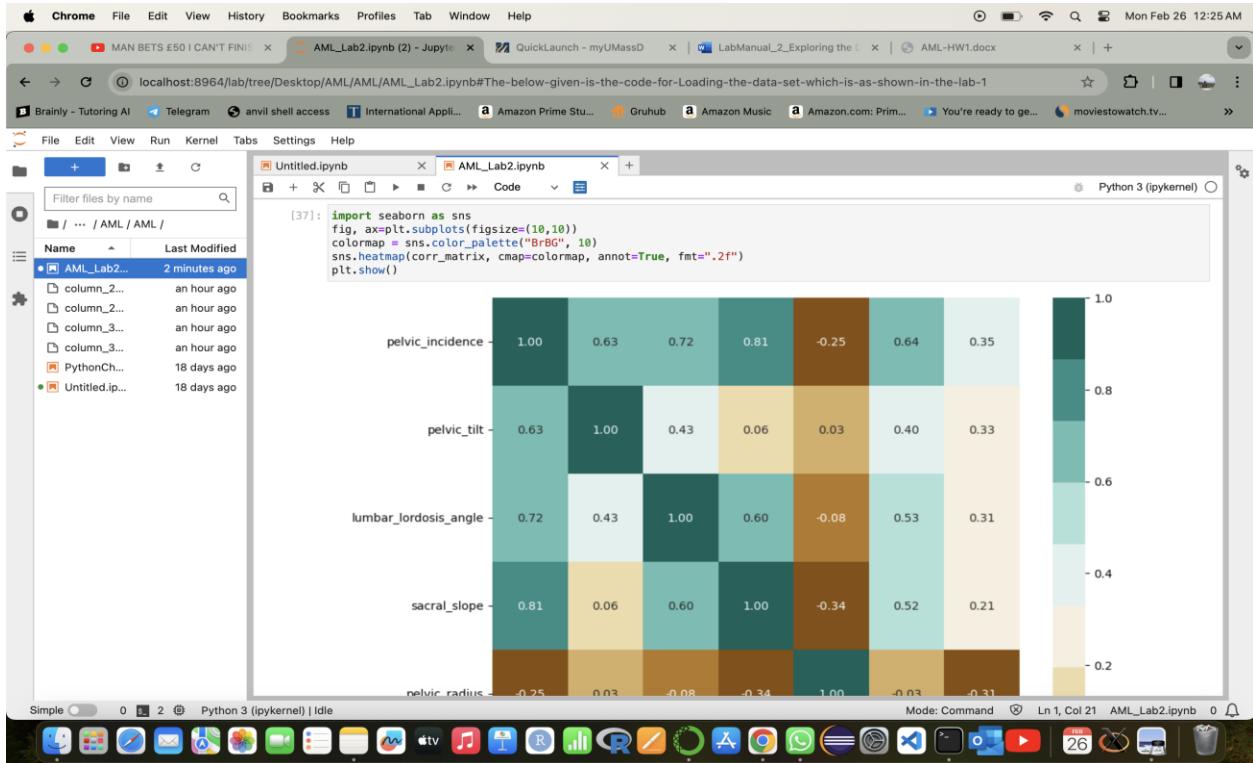
```
[68]: df.groupby('class').boxplot(fontsize=20,rot=90,figsize=(20,10),patch_artist=True)
[68]: 0      Axes(0,1,0,15;0.363636x0.75)
1      Axes(0.536364,0.15;0.363636x0.75)
dtype: object
```



This command will generate box plots for each numeric feature within each class category, allowing you to visually inspect differences in distributions, identify outliers, and get a sense of the spread and central tendency (median) of the data across different classes.

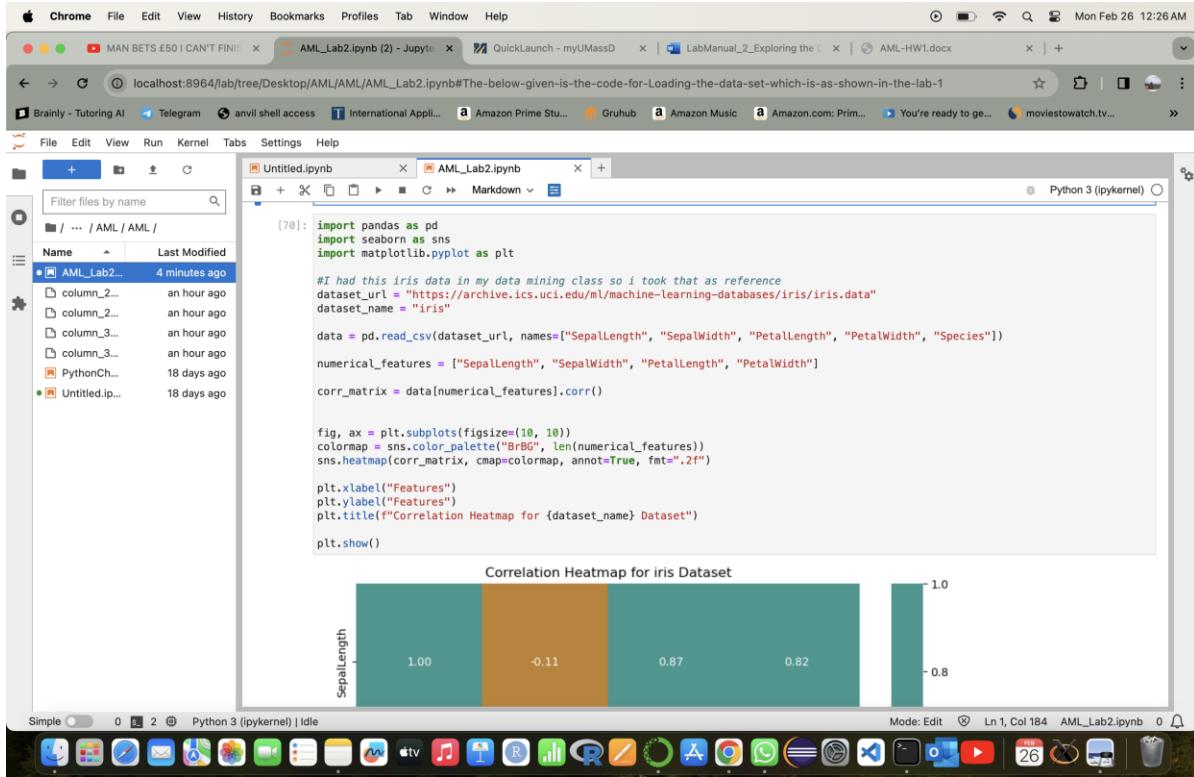


It uses pandas plotting capabilities to create a scatter matrix (also known as a pair plot) for the DataFrame df. This scatter matrix will include all the numerical columns in df.



The Seaborn library visualizes the correlation matrix of a dataset as a heatmap. This visualization is an effective way to quickly understand the relationships between variables in your dataset.

As given in the challenging task to "Find other data from the UCI Machine Learning Repository. Using the previous code for reference", The below is the code for that ML Repository



A screenshot of a Jupyter Notebook interface on a Mac OS X desktop. The notebook has two tabs: 'Untitled.ipynb' and 'AML_Lab2.ipynb'. The 'Untitled.ipynb' tab is active, showing Python code to load the Iris dataset and create a correlation heatmap. The code uses pandas to read the dataset, seaborn to create the heatmap, and matplotlib to display it. The resulting heatmap is titled 'Correlation Heatmap for iris Dataset' and shows the correlation coefficients between SepalLength, SepalWidth, PetalLength, and PetalWidth. The heatmap is color-coded with a 'BrBG' palette, where teal represents positive correlations and orange/brown represents negative correlations. The diagonal elements are all 1.0, and the off-diagonal elements show varying degrees of correlation, such as 0.87 for SepalLength and PetalLength, and -0.11 for SepalLength and SepalWidth.

```
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

# I had this iris data in my data mining class so i took that as reference
dataset_url = "https://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data"
dataset_name = "iris"

data = pd.read_csv(dataset_url, names=["SepalLength", "SepalWidth", "PetalLength", "PetalWidth", "Species"])

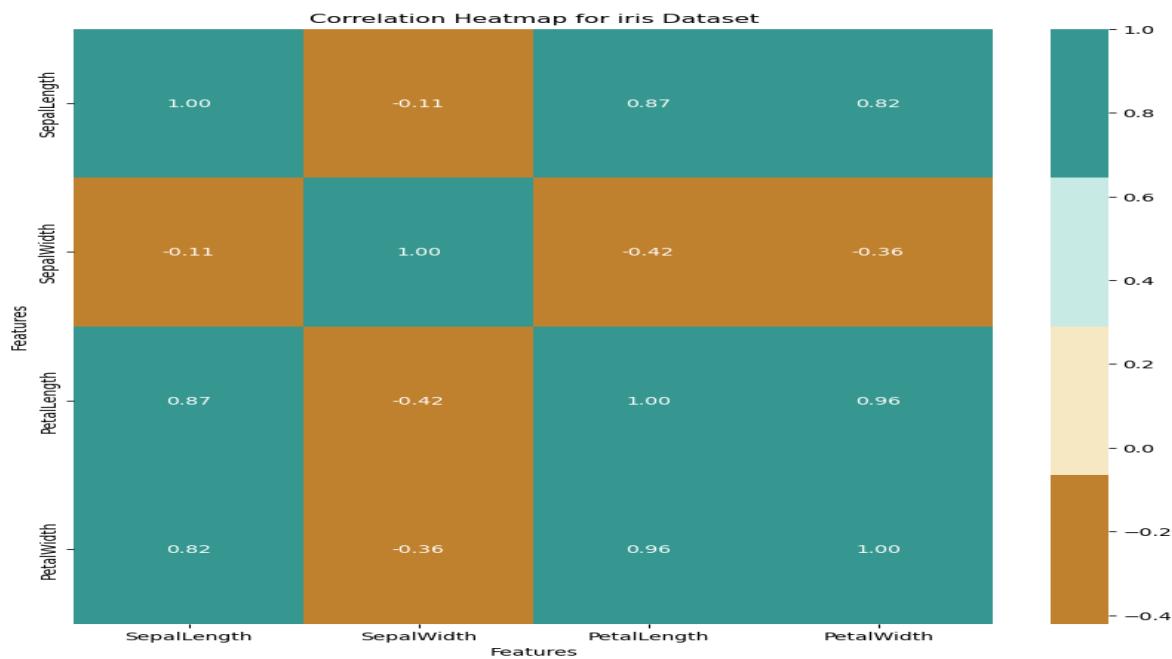
numerical_features = ["SepalLength", "SepalWidth", "PetalLength", "PetalWidth"]

corr_matrix = data[numerical_features].corr()

fig, ax = plt.subplots(figsize=(10, 10))
cmap = sns.color_palette("BrBG", len(numerical_features))
sns.heatmap(corr_matrix, cmap=cmap, annot=True, fmt=".2f")

plt.xlabel("Features")
plt.ylabel("Features")
plt.title(f"Correlation Heatmap for {dataset_name} Dataset")

plt.show()
```



It calculates the correlation matrix for the numerical features of the dataset and visualizes this matrix as a heatmap using Seaborn. This visualization helps in understanding the relationships between different numerical features of the Iris dataset.

Conclusion:

We have learned more about our dataset after doing this data exploration process:

- We'll learn about the dataset's size, structure, and features count
- Target variable analysis aids in understanding its distribution
- Class names may need numerical conversion for modeling.
- Data visualization aids in pattern comprehension.
- Scatter plots and heatmaps reveal numerical and categorical relationships.
- Outlier exploration helps assess their impact on analysis.
- Density plots and outlier detection unveil data patterns and possible anomalies.
- Data distribution insights come from numerical features' statistics using df.describe().
- Knowing the data types of columns guides our preprocessing steps.

By thoroughly investigating the data, we can make informed choices regarding data preprocessing, selecting the most pertinent characteristics, and developing modeling techniques. This comprehensive approach helps in creating more successful predictive models and extracting significant insights from the data.