

NAME : Maitri Vora
BRANCH : EXTC
UID : 2019120068
BATCH : B
EXPERIMENT : 1

AIM : To perform EDA such as number of data samples, number of features, number of classes, number of data samples per class, removing missing values, conversion to numbers, using the seaborn library to plot different graphs.

PROBLEM STATEMENT : Choose any one of the given datasets, apply EDA on it and write a detailed inference.

CODE & OUTPUT:

https://colab.research.google.com/drive/1YuyMm-eTKYxfI05_UkUUT8Tdt5l3Qw-?usp=sharing

CONCLUSION :

After learning about these exploratory data analysis techniques and also implementing them, I was able to conclude that :

- In the process of EDA, we first start by removing the redundant variables, that is, the columns that we know will not contribute to the model building and learning.
- The next step can comprise selecting the variable, which means studying every column and the kind of data that they possess, and also looking for ways to remove the null values.
- I studied the distribution of the features like mean, min, and max from the box plot which helped me determine whether the data contains outliers, and then eventually look for ways to remove them. In my case, I studied the value distribution using the `value_counts()` method, which showed me the data points after which the frequency becomes comparatively negligible and hence can be termed outlying. I used the Interquartile range method for removing the outliers.
- Finally, after data cleaning and preprocessing, I analysed the relationships between different variables and plotted various graphs and visualizations to understand their correlation. I used the seaborn library to plot graphs like heatmaps (it provides a coloured distribution that signifies a correlation between all the numerical data in the dataset). The observation has been mentioned in the colab notebook itself. After that, I plotted scatter plots, which showed a precise relationship between any two variables.