

EX.NO: 1	Study of Machine Learning Tools and Working with Dataset
DATE: 25.03.2022	

AIM:

To study and implement machine learning tools and working with datasets.

- 1. What is a dataset? Also create a sample dataset as csv file. Check the properties of a csv file.**

A dataset in machine learning is, quite simply, a collection of data pieces that can be treated by a computer as a single unit for analytic and prediction purposes. This means that the data collected should be made uniform and understandable for a machine that doesn't see data the same way as humans do

```
import csv
Head = ['football', 'volleyball', 'cricket', 'basketball']
Data = [ ['20EUA1005', '20EUA1010', '20EUA1011', '20EUA1006'],
        ['20EUA1016', '20EUA1003', '20EUA1007', '20EUA1022'],
        ['20EUA10020', '20EUA1018', '20EUA1025', '20EUA1030'],
        ['20EUA1026', '20EUA1015', '20EUA1033', '20EUA1037'],
        ['20EUA1044', '20EUA1040', '20EUA1029', '20EUA1035']]

with open('Score.csv', 'w', encoding='UTF8') as f:
    writer = csv.writer(f)
    writer.writerow(Head)
    writer.writerows(Data)

with open('Score.csv', 'rt') as f:
    data = csv.reader(f)
    for row in data:
        print(row)
```

THE PROPERTIES OF A CSV FILE:

A comma-separated values (CSV) file is a delimited text file that uses a comma to separate values. Each line of the file is a data record. Each record consists of one or more fields, separated by commas. The use of the comma as a field separator is the source of the name for this file format.

2. Perform a self-study on the python Packages used in Machine Learning.

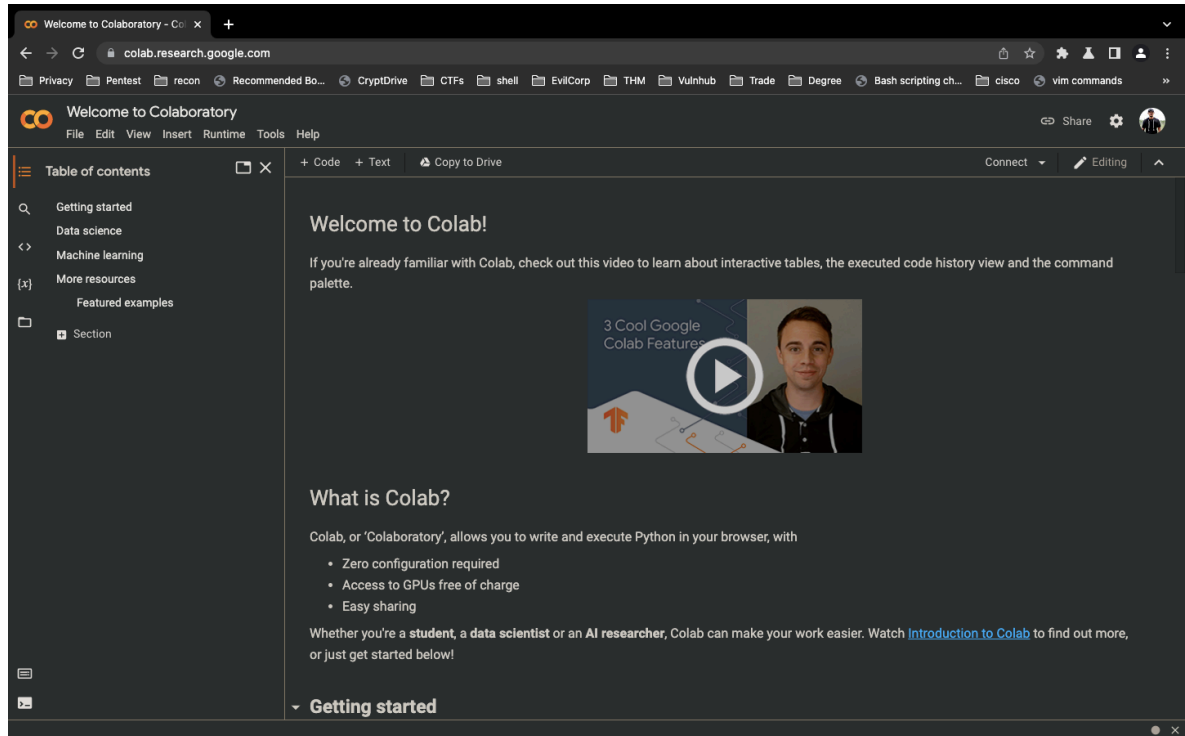
- a. Tensorflow
 - i. Tensorflow is a open source library which is used for Machine learning and Artificial intelligence which is used in python and java programming languages. It is focused on deep neural networks.
- b. Natural Language ToolKit (NLTK)
 - i. NLTK is the widely used library for Text Classification and Natural Language Processing.
- c. Sci-kit
 - i. Scikit-learn is mostly focused on various data modeling concepts like regression, classification, clustering, model selections
- d. Keras
 - i. Keras provides a Python interface of Tensorflow Library especially focused on AI neural networks.
- e. Pytorch
 - i. The main focus of the library is only on developing and training deep learning models only.
- f. MLpack
 - i. The main emphasis while developing this library was on making it a fast, scalable, and easy-to-understand as well as an easy-to-use library so that even a coder new to programming can understand and use it without any problem
- g. OpenCv
 - i. OpenCV is an open-source platform dedicated to computer vision and image processing

3. Give due importance to the following and see what the following modules and functions do:

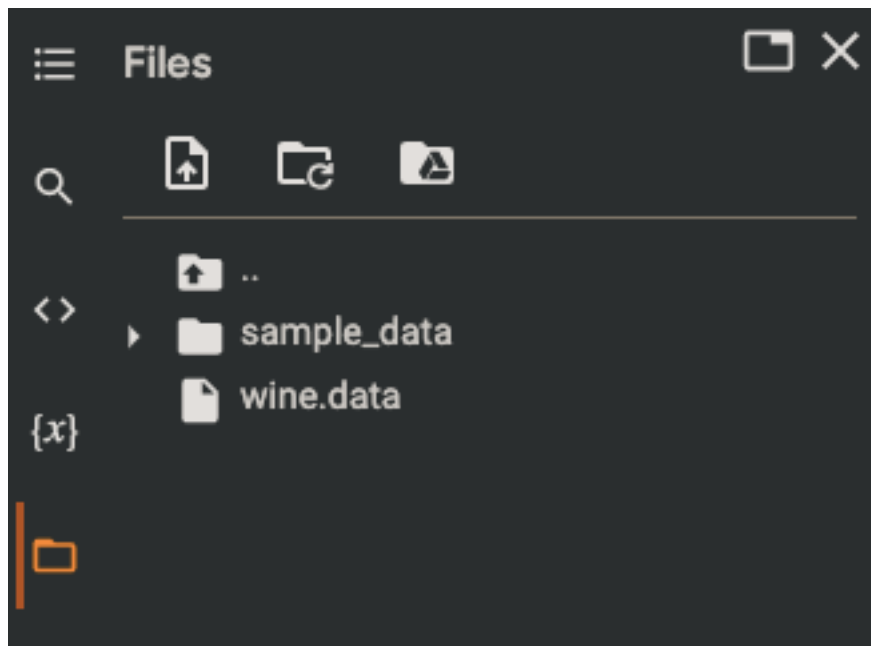
- a. Scipy
 - i. SciPy is a free and open-source Python library used for scientific computing and technical computing. SciPy contains modules for optimization, linear algebra, integration, interpolation, special functions, FFT, signal and image processing, ODE solvers and other tasks common in science and engineering.
- b. Numpy
 - i. NumPy is a library for the Python programming language, adding support for large, multi-dimensional arrays and matrices, along with a large collection of high-level mathematical functions to operate on these arrays.
- c. Pandas
 - i. pandas is a software library written for the Python programming language for data manipulation and analysis. In particular, it offers data structures and operations for manipulating numerical tables and time series
- d. Sklearn
 - i. Scikit-learn is a free software machine learning library for the Python programming language. It features various classification, regression and clustering algorithms including support-vector machines.
- e. Matplotlib
 - i. Matplotlib is a plotting library for the Python programming language and its numerical mathematics extension NumPy. It provides an object-oriented API for embedding plots into applications using general-purpose GUI toolkits like Tkinter, wxPython, Qt, or GTK.

4. Explore Google Colab and Load a dataset and see how it works.

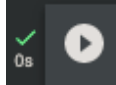
- a. This is how Google Colab looks like.



- b. Upload the csv file on the left menu.



- c. Write the python code in a code block and execute it using the play button on the left top corner of the code block cell.




```
import csv
Head = ['football', 'volleyball', 'cricket', 'basketball']
Data = [ ['20EUI005', '20EUI010', '20EUI011', '20EUI006'],
        ['20EUI016', '20EUI003', '20EUI007', '20EUI022'],
        ['20EUI0020', '20EUI018', '20EUI025', '20EUI030'],
        ['20EUI026', '20EUI015', '20EUI033', '20EUI037'],
        ['20EUI044', '20EUI040', '20EUI029', '20EUI035'] ]

with open('Score.csv', 'w', encoding='UTF8') as f:
    writer = csv.writer(f)
    writer.writerow(Head)
    writer.writerows(Data)

with open('Score.csv', 'rt') as f:
    data = csv.reader(f)
    for row in data:
        print(row)
```

- d. The output will be displayed below each cell block.



```
['football', 'volleyball', 'cricket', 'basketball']
['20EUI005', '20EUI010', '20EUI011', '20EUI006']
['20EUI016', '20EUI003', '20EUI007', '20EUI022']
['20EUI0020', '20EUI018', '20EUI025', '20EUI030']
['20EUI026', '20EUI015', '20EUI033', '20EUI037']
['20EUI044', '20EUI040', '20EUI029', '20EUI035']
```

5. Create a dataset with the following features:

Height(in cms)	Weight(in Kgs)
158	48
159	62
160	72
162	71
157	70
164	61
166	82
172	71
161	86
165	65

- a. Create the dataset in csv file and upload the file in Google Colab. Write the code and get the output.

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

data = pd.read_csv("example_dataset.csv")
data.head(10)
```

	Height (in cms)	Weight (in kgs)
0	158	48
1	159	62
2	160	72
3	162	71
4	157	70
5	164	61
6	166	82
7	172	71
8	161	86
9	165	65

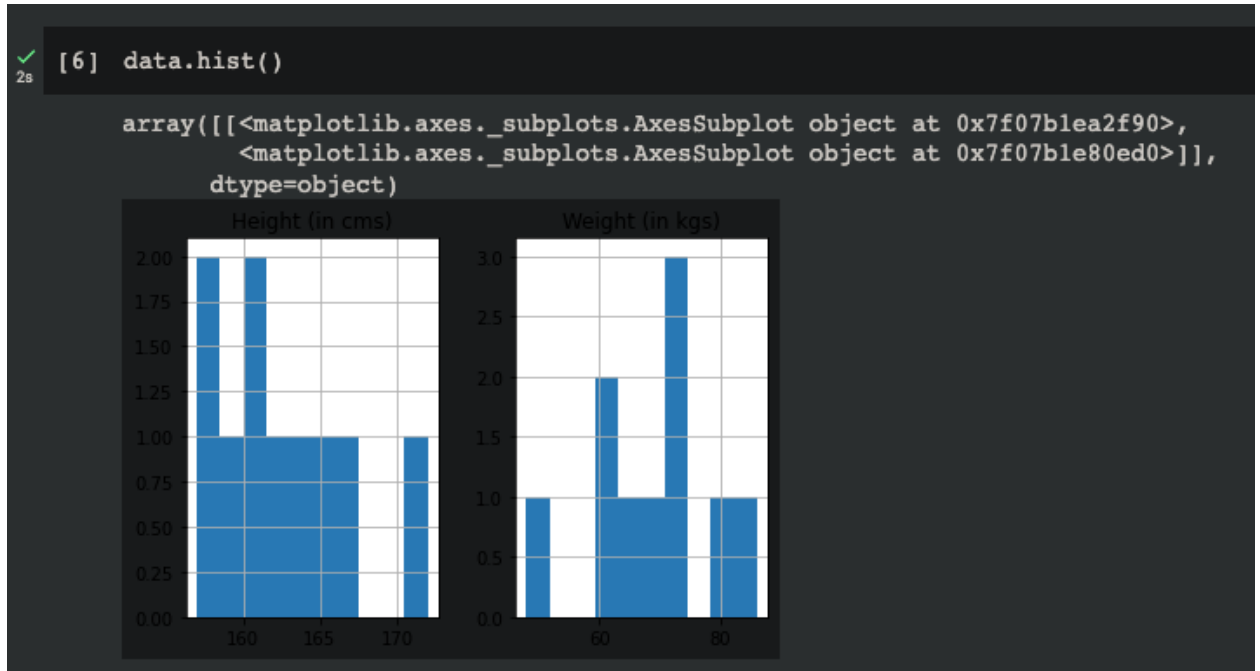
- b. Print some details about the csv file (dataset).

```
✓ 0s ▶ print(f"Shape of the data: {data.shape}")
      print(f"Type of the data: {type(data)}")
      print(f"Size of the data: {data.size}")

□ ▶ Shape of the data: (10, 2)
    Type of the data: <class 'pandas.core.frame.DataFrame'>
    Size of the data: 20
```

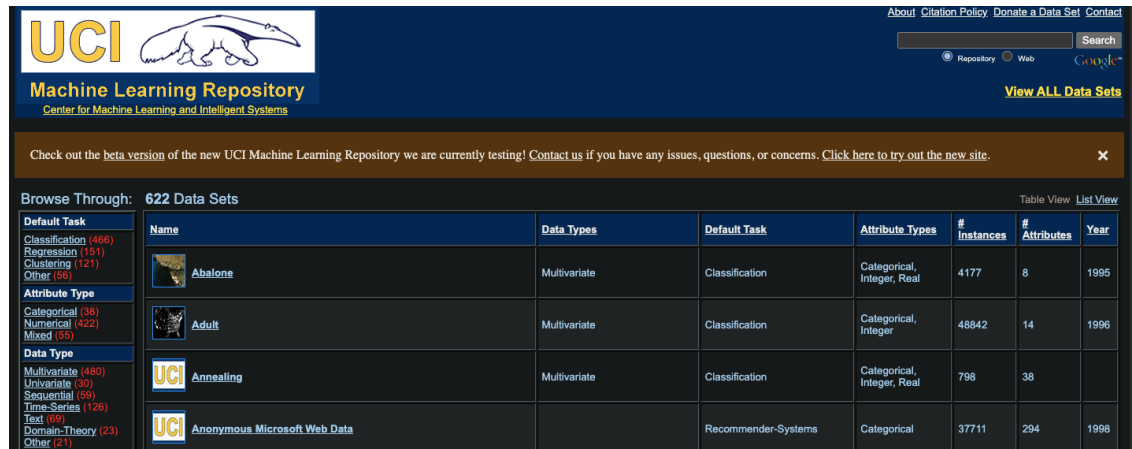
6. Identify the instances and plot a graph using these dimensions using matplotlib

Plotting histogram for the given csv file (dataset)



7. Explore more about importing datasets from UCI repositories and perform the same.

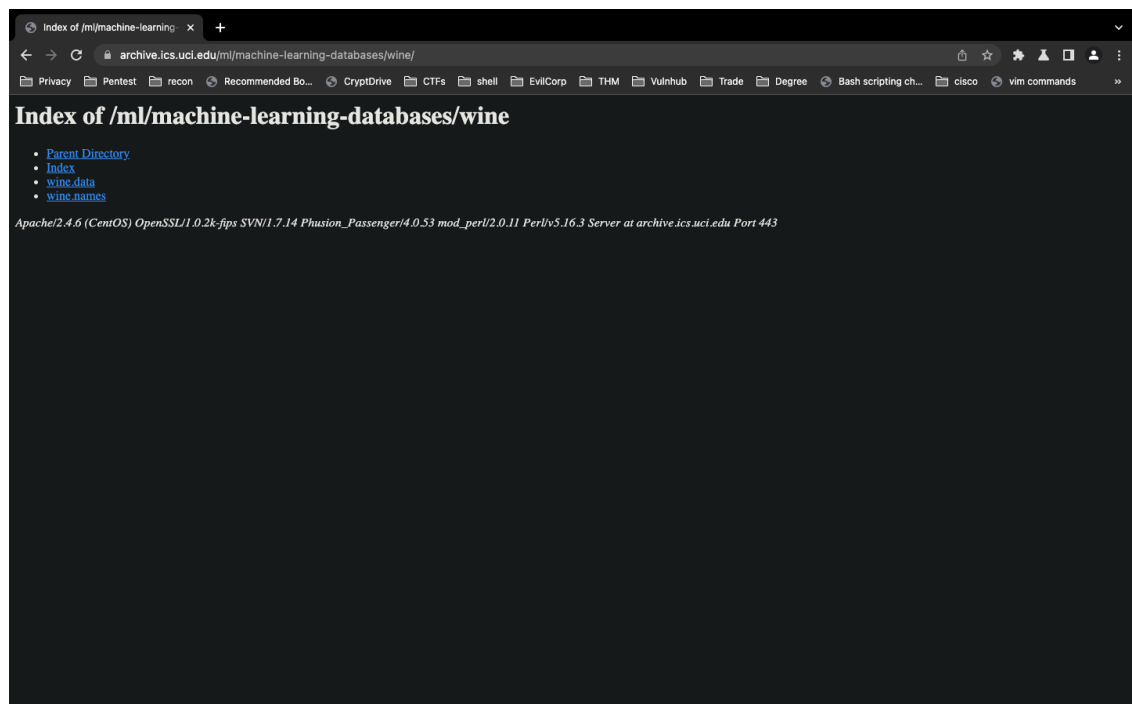
- Search for the required datasets from UCI repository.



The screenshot shows the UCI Machine Learning Repository website. The header includes the UCI logo, the text "Machine Learning Repository", and the subtitle "Center for Machine Learning and Intelligent Systems". There are links for "About", "Citation Policy", "Donate a Data Set", and "Contact". A search bar is present with a "Search" button and a "Google" logo. A banner below the header states: "Check out the beta version of the new UCI Machine Learning Repository we are currently testing! Contact us if you have any issues, questions, or concerns. Click here to try out the new site." Below the banner, there is a section titled "Browse Through: 622 Data Sets". On the left, there are filters for "Default Task" (Classification (466), Regression (151), Clustering (121), Other (56)), "Attribute Type" (Categorical (38), Numerical (422), Mixed (85)), and "Data Type" (Multivariate (480), Univariate (30), Sequential (59), Time Series (126), Text (69), Domain Theory (23), Other (21)). The main table lists datasets with columns: Name, Data Types, Default Task, Attribute Types, # Instances, # Attributes, and Year. The datasets listed are Abalone, Adult, Annealing, and Anonymous Microsoft Web Data.

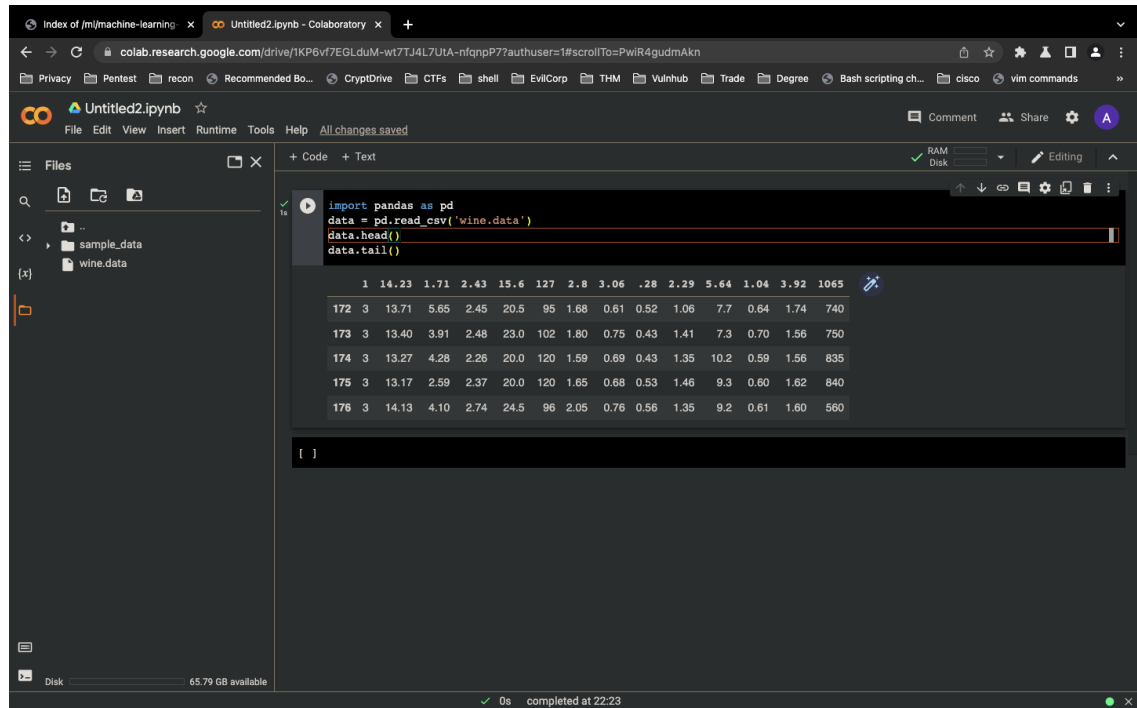
Name	Data Types	Default Task	Attribute Types	# Instances	# Attributes	Year
Abalone	Multivariate	Classification	Categorical, Integer, Real	4177	8	1995
Adult	Multivariate	Classification	Categorical, Integer	48842	14	1996
Annealing	Multivariate	Classification	Categorical, Integer, Real	798	38	
Anonymous Microsoft Web Data		Recommender-Systems	Categorical	37711	294	1998

- Download the required dataset.



The screenshot shows a web browser window with the address bar displaying "archive.ics.uci.edu/ml/machine-learning-databases/wine/". The page title is "Index of /ml/machine-learning-databases/wine". The page content includes a list of links: "Parent Directory", "Index", "wine.data", and "wine.names". Below the links, there is a footer line: "Apache/2.4.6 (CentOS) OpenSSL/1.0.2k-fips SVN/1.7.14 Phusion_Passenger/4.0.53 mod_perl/2.0.11 Perl/v5.16.3 Server at archive.ics.uci.edu Port 443".

c. Upload and run the python code to get the output.



The screenshot shows a Google Colaboratory notebook interface. The browser address bar indicates the notebook is located at `colab.research.google.com/drive/1KP6vt7EGLduM-wt7TJ4L7UIA-nfqnpP7?authuser=1#scrollTo=PwiR4gudmAkN`. The notebook is titled "Untitled2.ipynb" and has tabs for "Code" and "Text". The "Code" tab is active, showing the following Python code:

```
import pandas as pd
data = pd.read_csv('wine.data')
data.head()
data.tail()
```

The output of the code is displayed below the code cell. It shows the first five rows of the 'wine.data' CSV file. The data is presented in a table format with columns representing various wine attributes. The first row of data is:

	1	14.23	1.71	2.43	15.6	127	2.8	3.06	.28	2.29	5.64	1.04	3.92	1065
172	3	13.71	5.65	2.45	20.5	95	1.68	0.61	0.52	1.06	7.7	0.64	1.74	740
173	3	13.40	3.91	2.48	23.0	102	1.80	0.75	0.43	1.41	7.3	0.70	1.56	750
174	3	13.27	4.28	2.26	20.0	120	1.59	0.69	0.43	1.35	10.2	0.59	1.56	835
175	3	13.17	2.59	2.37	20.0	120	1.65	0.68	0.53	1.46	9.3	0.60	1.62	840
176	3	14.13	4.10	2.74	24.5	96	2.05	0.76	0.56	1.35	9.2	0.61	1.60	560

The notebook interface also shows a file explorer on the left with a folder named "sample_data" containing a file named "wine.data". The status bar at the bottom indicates that the code was completed at 22:23.

RESULT:

We have successfully studied and implemented machine learning tools and working with datasets.