## MACHINE LEARNING

GITHUB LINK: <a href="https://github.com/vijender6/vijender">https://github.com/vijender6/vijender</a>

RECORDING VIDEO LINK:

https://drive.google.com/file/d/1QfBjJF11W6inRNQWCPqluqrYbqFBxueg/view?usp=drive\_link

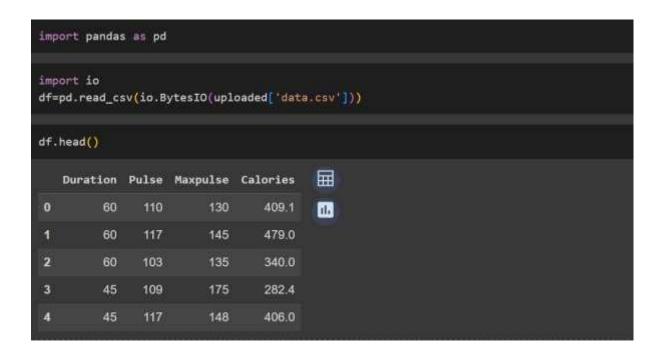
NAME: VIJENDER REDDY KOOTURU

700765220

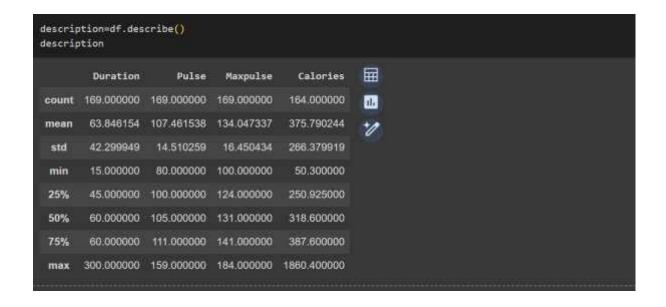
Q1. Read the provided CSV file 'data. csv'.

https://drive.google.com/file/d/1-Ir3AXK1A77A-qCDu5gGkAxv- nbmWlHO/view?usp=drive link

Here in the code, I have used  $read\_csv()$  and head() functions to read CSV file and to display the first 5 rows of the DataFrame respectively.



Q2. Here in the code, I have used the **describe()** function to generate descriptive statistics of a *DataFrame*. Then the description is printed.



Q3. Here in the code, I have used the *isnull*() function to identify null values in the *DataFrame*, and the *values*. *any*() function to check if any of the values are null.

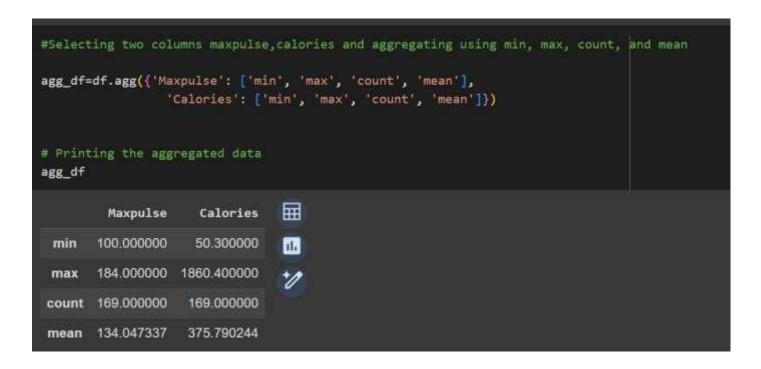
Then sum() function is called on this Boolean DataFrame and displayed a new DataFrame containing the sum of True values for each column in df.

```
[7]
    print( 'Are there any null values: ', df. isnull().values.any ())
    # Checking for null values in the Dataframe
    null_values = df. isnull(). sum()
    # Printing the number of null values for each column
    null_values
    Are there any null values:
    Duration
                 0
    Pulse
                 0
    Maxpulse
                0
    Calories
                 5
    dtype: int64
```

Q3(a). Replace the null values with the mean.

```
#Replacing null values with the mean of the respective column
mean_values = df.mean()
df. fillna(mean_values, inplace=True)
print('Are there any null values after replacing: ', df.isnull().values.any())
# Checking for null values in the Dataframe (shguld return all 0s)
null_values= df.isnull().sum()
# Printing the number of null values for each column
null values
Are there any null values after replacing: False
Duration
           0
Pulse
           0
           0
Maxpulse
Calories
           0
dtype: int64
```

Q4. I have used agg() function to perform the aggregation operations on the selected columns, with a dictionary as an argument, where the keys are the column names to be aggregated and the values are lists of aggregation functions to be applied to each column.



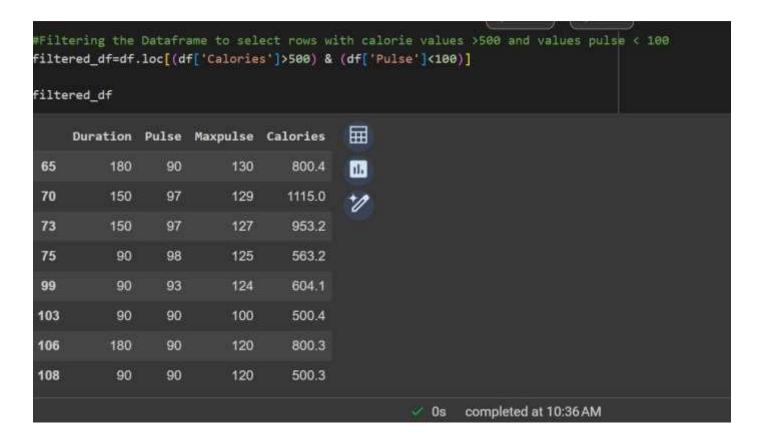
Q5. The loc[] function is used to select rows based on a boolean condition. The condition is specified inside the square brackets of the loc[] function using the 'Calories' column of df. Specifically, df['Calories'] >= 500 and df['Calories'] <= 1000 are two separate boolean conditions, which are combined using the & operator to specify the filter condition.

```
#Filtering the Dataframe to select rws with calorie values between 500 and 1000
filtered_df=df.loc[(df['Calories']>=500) & (df['Calories']<=1000)]
filtered_df</pre>
```

## OUTPUT:

51	80	123	146	643.1	11.
62	160	109	135	853.0	7
65	180	90	130	800.4	
66	150	105	135	873.4	
67	150	107	130	816.0	
72	90	100	127	700.0	
73	150	97	127	953.2	
75	90	98	125	563.2	
78	120	100	130	500.4	
83	120	100	130	500.0	
90	180	101	127	600.1	
99	90	93	124	604.1	
101	90	90	110	500.0	
102	90	90	100	500.0	

Q6. Filter the DataFrame to select the rows with calories values > 500 and pulse < 100.



Q7. Create a new " $df_{modified}$ " dataframe that contains all the columns from df except for "Maxpulse".

```
# Dropping the 'Maxpulse' column and creating a new Dataframe
df modified=df.drop(columns=['Maxpulse'])
#Printing the first 5 rows of new Dataframe
df_modified.head()
                                丽
             Pulse Calories
   Duration
0
          60
                110
                        409.1
                                 ılı
1
          60
                117
                        479.0
2
          60
                103
                        340.0
3
                        282.4
          45
                109
4
          45
                117
                        406.0
```

Q8. Delete the "Maxpulse" column from the main df dataframe.

```
# Dropping the 'Maxpulse' column from the orginal Dataframe
df.drop(columns=['Maxpulse'], inplace=True)
#Printing the first 5 rows of new Dataframe
df.head()
                              丽
   Duration Pulse Calories
0
         60
              110
                      409.1
                              ılı
1
         60 117
                      479.0
         60 103
                      340.0
2
3
         45 109
                      282.4
         45
                      406.0
4
              117
```

Q9. Convert the datatype of *Calories* column to *int* datatype.

```
#astype() function is used to convert the data type
df['Calories']=df['Calories'].astype('int64')
df.dtypes

Duration int64
Pulse int64
Calories int64
dtype: object
```