**Shourya Jindal**
**2020336**
**ML Assignment 1**

**SECTION A**

SHOURYA JINDAL
2020336
ML ASSIGNMENT 1

## SECTION A

(Q1) (a) No, a strong correlation b/w 2 variables with a 3rd variable. it does not necessarily apply that they will also display high degree of correlation with each other.

Explaination:

Correlation measures the statistical relation b/w 2 variables. But if the presence of 3rd variable that is related to both the other variables can affect the observed correlation b/w them. when 3rd variable is not taken into account it can create a misleading impression of direct reln b/w the 1st 2 variables.

Example:
$X$: Exercise (Variable representing how much a person is exercising)
$Y$: Diet (Variable representing how healthy diet a person is taking)
$Z$: Weight loss (Variable representing how much weight is lost by a person).
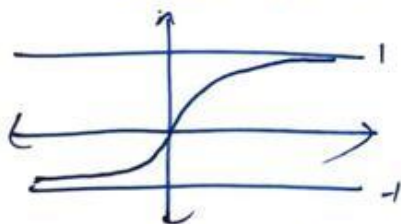
we, can see that generally
$X \propto Z$ (more a person exercises/workout, more he/she will lose weight)
$Y \propto Z$ (healthy diet $\Rightarrow$ more weight loss)

but we cannot say that $X$ and $Y$ are related. a person may Exercise more but still may not have good diet or vice-versa. Although, it may look like they are related since both leads to weight loss, but we have to consider other factors like "metabolism".
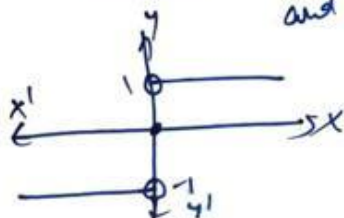
**(Q2)** Criterias for a mathematical func^n to be categorized as a logistic func^n are:

- Should have S-shaped curve, and must be bounded b/w a minima and maxima
- Domain should be $\mathbb{R}$ and Range must be $[0, 1]$ or bounded
- Should be continuous, differentiable
- We should be able to draw a decision boundary

- $Sinh(x)$ and $Cosh(x)$: They are not valid logistic func^n as they are not bounded

  Range of $Sinh x = (-\infty, \infty)$
  
  '' of $Cosh x = [1, \infty)$ } not bounded.

- $tanh(x)$: Yes it is valid (logistic func^n) as it satisfies all criterias.

  - It is sigmoidal i.e. S-shaped
  - Is bounded, continuous, differentiable
  - Domain $= \mathbb{R}$, Range $= [-1, 1]$ } bounded.



- $Signum(x)$: Not valid logistic func^n as it is not continuous and does not have S-shaped curve

**(Q3)** For very sparse datasets, leave one out cross Validation is beneficial.
This is because of the following reasons:

- It utilizes and maximises data usage for both training and validation
- It generally gives low ~~variance~~ bias estimate of the model. This is useful is sparse datasets as there is high degree of variability due to limited no. of samples

**How it is diff^n:**

- In this technique, we train the model on all the data pts except one, which is then used for testing. This process is repeated till all pts have been used for testing.
  Average performance is calculated for all the iterations.

- While in k-fold. we divide dataset into k-subsets and then uses 'k-1' folds for training and '1' fold for testing. This is repeated 'k' times.

| K-fold | Leave one out |
|---|---|
| • Requires K-iterations so, faster | • Req. n-iterations, hence requires more time and computes. |
| • balance b/w bias and variance depending on K. | • Generally low bias and high variance model, as it trains on all the data except one |

(Q4) Find coeff$^n$ of least square regression in slope-intercept form.

Let reg func$^n$ be: $\quad y = mx + c$

where m = slope and c = intercept
are the unknown coefficients.

Cost func$^n$ for least square regression would be:

$$J(y,n) = \frac{1}{n} \sum_{i=1}^{n} (y_i - y)^2 \qquad (\text{for } n - \text{data points})$$

$y_i \rightarrow$ actual value
$y \rightarrow$ predicted value
$n_i \rightarrow$ input

⊗ we need to minimize $J(y,n)$.

So, minimize $J(y,n) = \frac{1}{n} \sum_{i=1}^{n} (y_i - mx_i + c)^2$

So,

$$\frac{\delta J}{\delta m} = 0 \qquad \& \qquad \frac{\delta J}{\delta c} = 0$$

$\Rightarrow \frac{2}{n} \sum (y_i - mx_i + c)(-x_i) = 0$  $\quad \& \quad \frac{2}{n} \sum (y_i - mx_i + c) = 0$

$\Rightarrow \boxed{m = \dfrac{\sum x_i y_i - c \sum x_i}{\sum x_i^2}}$  $\quad \& \quad \boxed{m = \dfrac{\sum y_i - nc}{\sum x_i}}$

on equating we set $\boxed{m = \dfrac{n \sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2}}$

$m = \cancel{\sum x_i} \quad c = c \sum y_i + A$

and $\boxed{c = \dfrac{\sum y_i - m \sum x_i}{n} = \dfrac{\sum y_i}{n} - \left(\dfrac{\sum x_i}{n}\right)\left(\dfrac{n \sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2}\right)}$

(Q5) Ans: (a) $\alpha, \beta, \sigma$
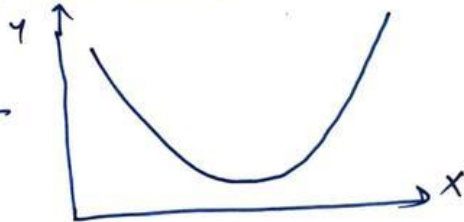
$y = \alpha + \beta x + \epsilon \quad \epsilon \sim N(0, \sigma)$

Here '$\alpha, \beta$' are the coefficients/weights/parameters of
simple linear regression model and
'$\sigma$' is the standard deviation of the $\epsilon$ variable
within which is the Noise parameter which
follows $N(0, \sigma)$ dist".

(Q6) Ans: (d) $y = \alpha + \beta_1 x + \beta_2 x^2 + \epsilon \quad \beta_2 > 0$

Reason: $X = [20, 30, 50, 60, 80, 90]$
$Y = [125, 110, 95, 90, 110, 130]$
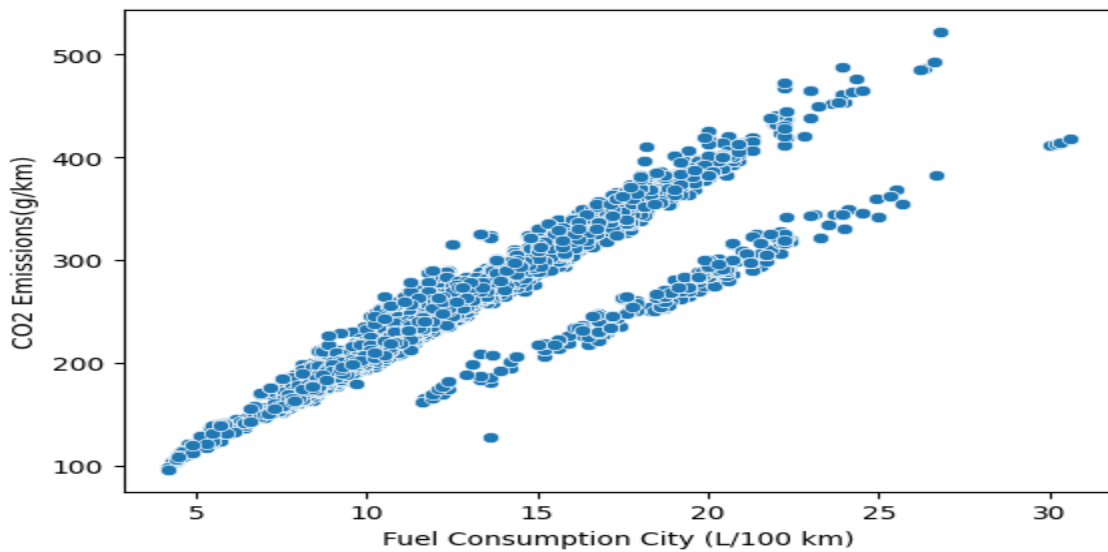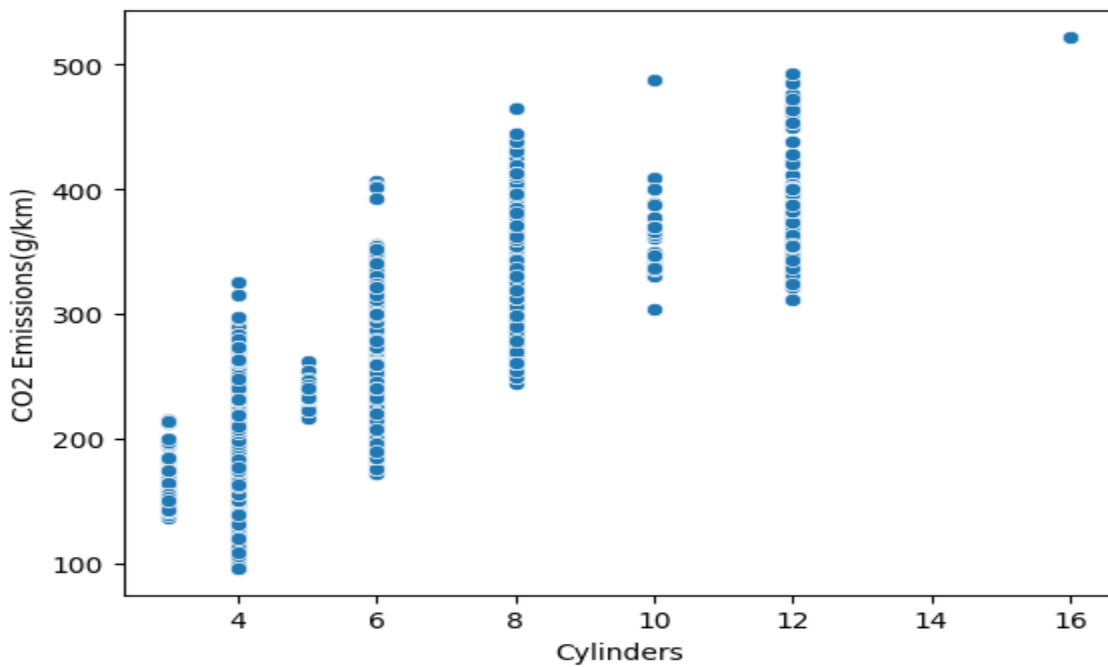
If we plot this graph:

Rough ← 
graph

as we can see it is
not linear but quadratic
hence (a) (b) are ruled out
also, this parabolic/quadratic graph reaches
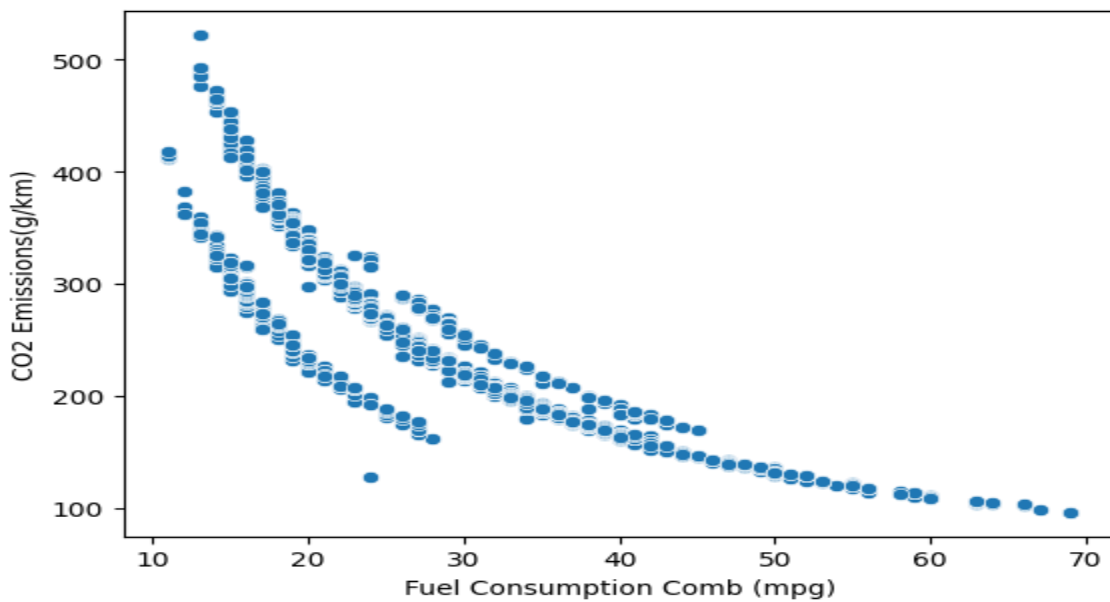minima and is upward zero facing
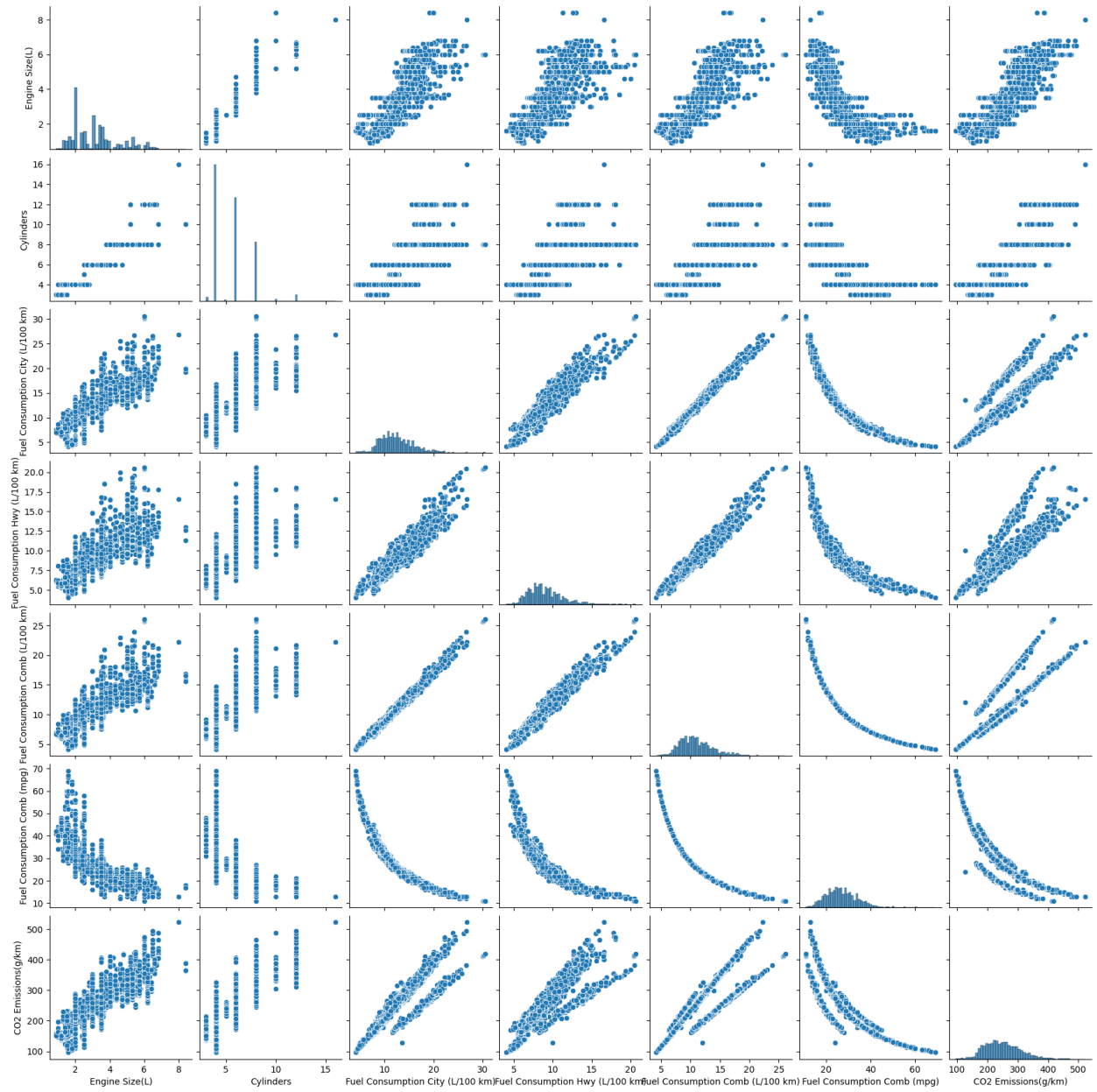thus. $\beta_2 > 0$.

# SECTION C

**PART A: DATA VISUALIZTION**

I visualized the data by creating several graphs and plots.
Some of which are:

Correlation Heatmap

FUEL TYPE WITH PERCENTAGES

Rest are in the other folder and in the code itself.
Some insights for the data are:

- Cylinders are directly proportional to CO2 Emissions
- Engine Size is also directly proportional to CO2 Emissions.
- From heatmap we can see that cylinders and engine size has a very high correlation.
- Also Fuel Consumption (city, hwy etc.) are also directly proportional to CO2 emissions
- Almost 50% of the vehicles are of Fuel Type - 'X', and 43% of type 'Z' and rest have very few numbers.
- From the box plot of Vehicle class vs CO2 emissions, we infer that vehicle class do not affects CO2 emissions very much.
- Also 'make', 'model' have very low correlation with Co2 emissions.

**Part B: TSNE**
I used TSNE algorithm (from sklearn library) to reduce data dimensions to 2 and plotted the resulting data as a scatter plot.

TSNE SCATTER PLOT

The scatter plot shows that most of the clusters are not independent. The features forms clusters which intermixing on TSNE scatter plot and thus, are not seperable and are most of them are interrelated.

## PART C: LABEL ENCODING
Used LabelEncoding from sklearn to LabelEncode the categorical data : `'Make', 'Model', 'Vehicle Class', 'Transmission', 'Fuel Type']`

Then after preprocessing, I split the data into train:test = 80:20, then applied Linear Regression on it, also using sklearn.

The Metric and Performance Results are as follows:

|  | TRAINING DATA | TESTING DATA |
|---|---|---|
| MSE | 295.177056 | 258.138270 |
| RMSE | 17.180718 | 16.066682 |
| R2 | 0.913177 | 0.926511 |
| Adjusted R2 | 0.913015 | 0.926373 |
| MAE | 11.238194 | 10.470382 |

We get a good R2 score of 0.91 and very less rmse error.

**PART D: PCA on label encoded**

Results are as follows:

```
PCA WITH NO OF COMPONENTS = 2
                TRAINING DATA   TESTING DATA
MSE             2776.707846     2554.442701
RMSE              52.694476       50.541495
R2                 0.200487        0.207222
Adjusted R2        0.200216        0.206954
MAE               41.293940       39.795481

PCA WITH NO OF COMPONENTS = 4
                TRAINING DATA   TESTING DATA
MSE              457.237671      436.144308
RMSE              21.383116       20.884068
R2                 0.867724        0.867347
Adjusted R2        0.867634        0.867257
MAE               13.783732       13.475276

PCA WITH NO OF COMPONENTS = 6
                TRAINING DATA   TESTING DATA
MSE              369.059370      393.693058
RMSE              19.210918       19.841700
R2                 0.891220        0.888947
Adjusted R2        0.891109        0.888834
MAE               11.030404       11.427482

PCA WITH NO OF COMPONENTS = 8
                TRAINING DATA   TESTING DATA
MSE              287.753851      302.565302
RMSE              16.963309       17.394404
R2                 0.916848        0.907566
Adjusted R2        0.916735        0.907441
MAE               11.119796       11.333621

PCA WITH NO OF COMPONENTS = 10
                TRAINING DATA   TESTING DATA
MSE              291.276442      274.122151
RMSE              17.066823       16.556635
R2                 0.915465        0.917736
Adjusted R2        0.915322        0.917597
MAE               11.155951       11.018914
```

Used PCA from sklearn to implement this. We observe that as no of components increases error decreases and accuracy increases. The R2 score increased from 0.20 (when no of components = 2) to 0.91(when no of components = 10), which is very significant improvement.

**PART E: One-Hot Encoding**

I one hot coded the original data  on categorical data using pd.get_dummies.
Since the no. of distinct values of categorical data was very large, no of columns increases to almost 2150 thus greatly increasing the size of the data.
Then we did Linear Regression using sklearn.

Performance analysis:

|  | TRAINING DATA | TESTING DATA |
|---|---|---|
| MSE | 8.570359 | 2.888609e+20 |
| RMSE | 2.927518 | 1.699591e+10 |
| R2 | 0.997464 | -8.030857e+16 |
| Adjusted R2 | 0.996014 | -1.262328e+17 |
| MAE | 1.893353 | 3.508334e+09 |

On comparing the these results with part c, we observe that one hot encoding performs much better than label encoding on training data but performs very poorly on testing data as compared to part c.
This is because one hot encoding led the model to get overfit. As a result bais decreased but variance increased very much.

**PART F: PCA on One-hot Encoded:**

Performance analysis:

```
PCA WITH NO OF COMPONENTS = 2
             TRAINING DATA   TESTING DATA
MSE              376.272097     383.239369
RMSE              19.397734      19.576500
R2                 0.889951       0.888564
Adjusted R2        0.889914       0.888527
MAE               11.006305      11.030149


PCA WITH NO OF COMPONENTS = 4
             TRAINING DATA   TESTING DATA
MSE              336.824469     318.892616
RMSE              18.352778      17.857565
R2                 0.902115       0.904769
Adjusted R2        0.902049       0.904704
MAE               11.591412      11.440505


PCA WITH NO OF COMPONENTS = 6
             TRAINING DATA   TESTING DATA
MSE              325.036415     320.624306
RMSE              18.028766      17.905985
R2                 0.905821       0.903123
Adjusted R2        0.905726       0.903025
MAE               11.423221      11.482435


PCA WITH NO OF COMPONENTS = 8
             TRAINING DATA   TESTING DATA
MSE              323.992192     322.473531
RMSE              17.999783      17.957548
R2                 0.904944       0.907380
Adjusted R2        0.904815       0.907254
MAE               11.387942      11.514365


PCA WITH NO OF COMPONENTS = 10
             TRAINING DATA   TESTING DATA
MSE              324.861023     318.542335
RMSE              18.023901      17.847754
R2                 0.905793       0.904132
Adjusted R2        0.905634       0.903969
MAE               11.405065      11.274745
```

As the no of components increases performance slightly get improved.

When we compare training and testing metrics, we can see that they are almost the same and there is not much difference. This is because the data is quite large, also complexity of the model is accurate. Also, shows that methods, models and metrics and evaluations used are good enough.

## PART G: L1 and L2 Regularization

L1 - Lasso and L2 - Ridge

Results For L1:

```
                TESTING DATA
MSE                263.356256
RMSE                16.228255
R2                   0.923112
Adjusted R2          0.922969
MAE                 10.519545
```

Results For L2:

```
                TESTING DATA
MSE                266.076842
RMSE                16.311862
R2                   0.922318
Adjusted R2          0.922173
MAE                 10.640248
```

Used label encoded data from part c for this part. Then used Lasso() and Ridge() from sklearn to perform this.
On comparing, we see that performance for both the methods are almost similar and also does not improve the results from part c.
This shows that none of this is useful and data is not overfitted.

**PART H: SGDRegressor:**

Results:

```
               TRAINING DATA   TESTING DATA
 MSE             4.209948e+28   4.184357e+28
 RMSE            2.051816e+14   2.045570e+14
 R2             -1.236362e+25  -1.198132e+25
 Adjusted R2    -1.238669e+25  -1.200367e+25
 MAE             1.782327e+14   1.764734e+14
```

Used SGDRegressor from sklearn on label encoded data from part c.
We see that performance is very very poor. Error is very high and R2 scores and very less.
This is because SGD is not optimal many times. To save time and computations it does not
reaches optimal performance.