

## HW 3

This assignment covers several aspects of Linear Regression. **DO NOT ERASE MARKDOWN CELLS AND INSTRUCTIONS IN YOUR HW submission**

- **Q** - QUESTION
- **A** - Where to input your answer

## Instructions

Keep the following in mind for all notebooks you develop:

- Structure your notebook.
- Use headings with meaningful levels in Markdown cells, and explain the questions each piece of code is to answer or the reason it is there.
- Make sure your notebook can always be rerun from top to bottom.
- Please start working on this assignment as soon as possible. If you are a beginner in Python this might take a long time. One of the objectives of this assignment is to help you learn python and scikit-learn package.
- Follow [README.md \(README.md\)](#) for homework submission instructions

## Tutorials

- [scikit-learn linear model \(https://scikit-learn.org/stable/modules/linear\\_model.html\)](https://scikit-learn.org/stable/modules/linear_model.html)
- [train-test-split \(https://towardsdatascience.com/train-test-split-and-cross-validation-in-python-80b61beca4b6\)](https://towardsdatascience.com/train-test-split-and-cross-validation-in-python-80b61beca4b6)
- [least squares fitting \(https://python4mpia.github.io/fitting\\_data/least-squares-fitting.html\)](https://python4mpia.github.io/fitting_data/least-squares-fitting.html)
- [Linear Regression \(https://scikit-learn.org/stable/modules/generated/sklearn.linear\\_model.LinearRegression.html\)](https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LinearRegression.html)
- [Seaborn \(https://seaborn.pydata.org/api.html\)](https://seaborn.pydata.org/api.html)

## REGRESSION TASK USING SKLEARN

In jupyter notebook environment, commands starting with the symbol % are magic commands or magic functions. `%%timeit` is one of such function. It basically gives you the speed of execution of certain statement or blocks of codes.

```
In [48]: ▶ import pandas as pd
import numpy as np
import seaborn as sns
```

Type *Markdown* and LaTeX:  $\alpha^2$

**\*Data \*** Get the exploratory data and the following files:

<https://archive.ics.uci.edu/ml/machine-learning-databases/auto-mpg/auto-mpg.data>  
<https://archive.ics.uci.edu/ml/machine-learning-databases/auto-mpg/auto-mpg.data>  
<https://archive.ics.uci.edu/ml/machine-learning-databases/auto-mpg/auto-mpg.names>  
<https://archive.ics.uci.edu/ml/machine-learning-databases/auto-mpg/auto-mpg.names>

**or** Use from our 2022Spring/data repository folder

- Link should automatically download the data
- copy them in your HW folder
- If you are using command line: `>> wget https://archive.ics.uci.edu/ml/machine-learning-databases/auto-mpg/auto-mpg.data`
  - If wget is not working
    - download it from [link \(https://eternallybored.org/misc/wget/\)](https://eternallybored.org/misc/wget/)
    - follow [steps \(https://stackoverflow.com/questions/29113456/wget-not-recognized-as-internal-or-external-command\)](https://stackoverflow.com/questions/29113456/wget-not-recognized-as-internal-or-external-command)

**Q1** Read the data using pandas, and replace the ??? in the code cell below to accomplish this task. Note that auto-mpg.data does not have the column headers. use auto-mpg.names file to provide column names to the dataframe.

**A1**

```
In [97]: # Replace ??? with code in the code cell below
column_names = ["mpg", "cylinders", "displacement", "horsepower", "weight", "acceleration", "model year", "origin", "car name"]
df = pd.read_csv('auto-mpg.data', names=column_names, na_values = "?",
```

```
In [98]: # View head of the data to confirm the correctness of your answer
```

Out[98]:

	mpg	cylinders	displacement	horsepower	weight	acceleration	model year	origin	car name
0	18.0	8	307.0	130.0	3504.0	12.0	70	1	NaN
1	15.0	8	350.0	165.0	3693.0	11.5	70	1	NaN
2	18.0	8	318.0	150.0	3436.0	11.0	70	1	NaN
3	16.0	8	304.0	150.0	3433.0	12.0	70	1	NaN
4	17.0	8	302.0	140.0	3449.0	10.5	70	1	NaN

## Data cleaning and manipulation

Use

**Q2** Data cleaning and manipulation:

1. use `pandas.info()` method to find columns with large number of NaN values
2. remove the column with NaN values

3. Check if there are still NaN values in the dataframe using `isna()` method

**A2** Replace ??? with code in the code cell below

```
In [99]: #1. use pandas.info() method to find columns with large number of NaN
df.info()

#2. remove the column with NaN values - replace ??? with code
df.drop(columns = ['car name'], inplace = True)
# Print head
df.head()

#3. Check if there are still NaN values in the dataframe using ``isna
df.isna().sum().sum()
# drop if any left or replace Nan values
df.dropna(inplace = True)
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 398 entries, 0 to 397
Data columns (total 9 columns):
#   Column                Non-Null Count  Dtype
---  -
0   mpg                   398 non-null   float64
1   cylinders              398 non-null   int64
2   displacement           398 non-null   float64
3   horsepower             392 non-null   float64
4   weight                 398 non-null   float64
5   acceleration           398 non-null   float64
6   model year            398 non-null   int64
7   origin                 398 non-null   int64
8   car name              0 non-null     float64
dtypes: float64(6), int64(3)
memory usage: 28.1 KB
```

```
In [100]: #Print Tail
df.head(15)
```

Out[100]:

	mpg	cylinders	displacement	horsepower	weight	acceleration	model year	origin
<b>383</b>	38.0	4	91.0	67.0	1965.0	15.0	82	3
<b>384</b>	32.0	4	91.0	67.0	1965.0	15.7	82	3
<b>385</b>	38.0	4	91.0	67.0	1995.0	16.2	82	3
<b>386</b>	25.0	6	181.0	110.0	2945.0	16.4	82	1
<b>387</b>	38.0	6	262.0	85.0	3015.0	17.0	82	1
<b>388</b>	26.0	4	156.0	92.0	2585.0	14.5	82	1
<b>389</b>	22.0	6	232.0	112.0	2835.0	14.7	82	1
<b>390</b>	32.0	4	144.0	96.0	2665.0	13.9	82	3
<b>391</b>	36.0	4	135.0	84.0	2370.0	13.0	82	1
<b>392</b>	27.0	4	151.0	90.0	2950.0	17.3	82	1

	mpg	cylinders	displacement	horsepower	weight	acceleration	model year	origin
393	27.0	4	140.0	86.0	2790.0	15.6	82	1
394	44.0	4	97.0	52.0	2130.0	24.6	82	2
395	32.0	4	135.0	84.0	2295.0	11.6	82	1

**Q3:**

1. Convert following columns 'cylinders', 'year', 'origin' to dummy variable using pandas `get_dummies()` function
2. Do data normalization on real value/continous columns
  - The formula for normalization is:  $(\text{Col\_value} - \text{Mean of the col}) / \text{Standard Deviation of the col}$

**A3** Replace ??? with code in the code cell below

```
In [101]: ▶ # 1. Convert following columns 'cylinders', 'year', 'origin' to dummy
cols = ['cylinders', 'model year', 'origin']
df_dummies = pd.get_dummies(df, columns = cols, prefix = ['cylinders',

#show the head
df_dummies.head()

# 2. Do data normalization on real value/continous columns
realcols = ['mpg', 'displacement', 'horsepower', 'weight', 'acceleration']

for col in realcols:
    mean = df[col].mean()
    std = df[col].std()
    df[col] = (df[col] - mean) / std
```

## Regression Task

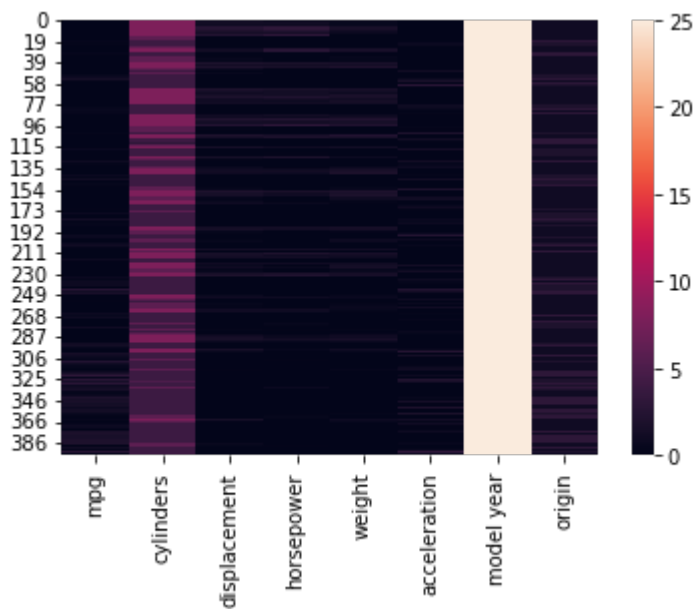
Given all the information we will try to predict mpg - miles per gallon. The First step toward predicting the mpg from the dataset is to find the correlation between the columns/features of the dataset.

**Q4**

1. Use heatmap chart from seaborn library to findout the correlation between the columns.
2. Which of the columns is mostly related to mpg column and why?

```
In [102]: # A4 code goes below
```

```
Out[102]: <AxesSubplot:>
```



#### A4

Looking at the heatmap I would say acceleration and mpg have the closest relation just based on the line pattern and color they contain.

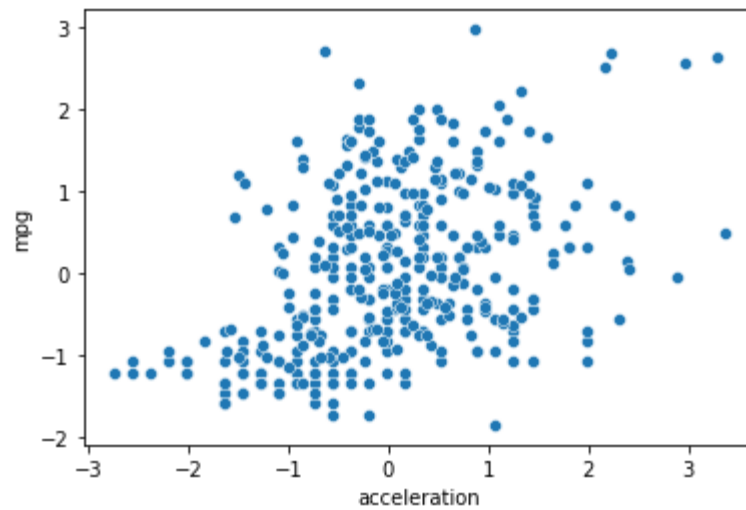
#### Q5

1. Draw a lineplot or scattered plot between mpg and your answer from the above cell.
2. Use pairplot from sns to plot our data frame df for better understanding of your selection
  - NOTE: 2. should inform 3.
3. Choose a set of columns which are significantly related to our goal.
  - Justify your answer using some explanation from the correlation chart above.
  - HINT: Refer to the newspaper radio TV example from the class slides

**A5** For 1. and 2. replace ??? with code in the code cell below.

```
In [103]: # 1. Draw a lineplot or scattered plot between mpg and your answer from
sns.scatterplot(data = df, x = 'acceleration', y = 'mpg')
# 2. Use pairplot from sns to plot our data frame df for better unders
sns.pairplot(df)
```

Out[103]: <seaborn.axisgrid.PairGrid at 0x1dd65c678b0>





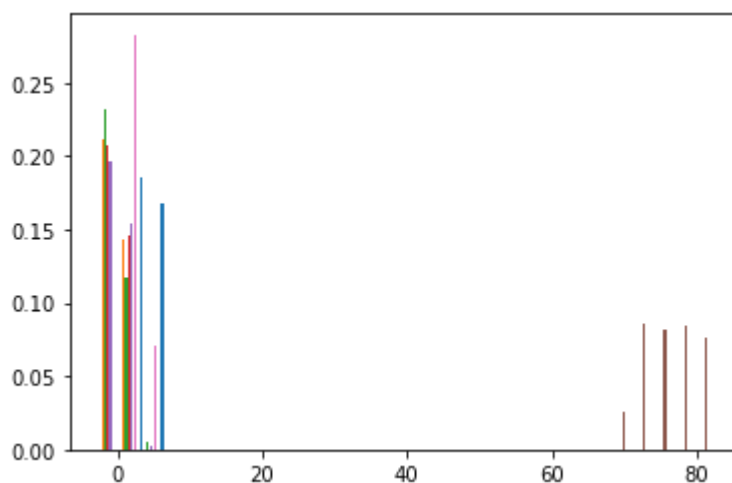
**A5** For 3., Looking into pairplot it would seem that there are two columns that are signficaly related to mpg. Those being weight and horsepower. My original answer doesn't seem to stand here as these two columns have a much tighter line. The closer the data points are to each other the more tightly connected the data is.

**Q6** Data Visualization:

1. Now, create a histogram which represents number items with per cylinder class

**A6** Replace ??? with code in the code cell below

```
In [95]: plt.hist(df, density = True, bins =30)
```



## Data Preparation

**Q7** Assign mpg column value to y and rest columns to x, remember x shouldn't have mpg

**A7** Replace ??? with code in the code cell below

```
In [104]: y = df.mpg.values
df.drop(columns=['mpg'],inplace = True)
```

**Q8** Use train\_test\_split to split the data set as train:test=(80%:20%) ratio.

**A8** Replace ??? with code in the code cell below

```
In [105]: from sklearn.model_selection import train_test_split

xtrain, xtest, ytrain, ytest = train_test_split(x,y,test_size= 0.2)
# View the shape of your data set
```

```
Out[105]: ((313, 7), (79, 7), (313,), (79,))
```

**Q9** Follow examples from references given in the top of this notebook

- Note: Use linear model to fit regression line and plot
- Our linear model will be of following type
- $Y = b + \text{coef}_0x_0 + \text{coef}_1x_1 + \text{coef}_2x_2 + \dots$

**A9:** Replace ??? with code in the code cell below

```
In [106]: > from sklearn import linear_model

reg = linear_model.LinearRegression()
reg.fit(xtrain,ytrain)

#Now view the coefficient use .coef_ and shape of .coef_
print(reg.coef_)
print(reg.intercept_)

[-0.07798248  0.26526127 -0.06065668 -0.70495678  0.02276379  0.09732
893
 0.1800173 ]
-7.2539802482502544
7
```

**Q10** Relates to the code in the cell below. Why the printed values the same?

```
In [107]: > # Now if you view
print(f'{reg.coef_.shape[0]},{xtrain.shape[1]}, ', f'are equal? {reg.c

7,7,  are equal? True
```

**A10** The numbers are print the same value due to coef\_ storing the fit method arrays of both x and y. Xtrain was one of the fit methods thus it's stored in coef\_

## Model Scoring

```
In [108]: > # Model Score
from sklearn import linear_model
reg = linear_model.LinearRegression()
reg.fit(xtrain, ytrain)
reg.score(xtest,ytest)

# Calculate the score on train and test sets
# Your code goes below
```

```
Out[108]: (0.8154769852171251, 0.8449374024375086)
```

**Q11** Each of the sklearn models have different model evaluations core value.

- LinearRegression [documentation \(https://scikit-learn.org/stable/modules/generated](https://scikit-learn.org/stable/modules/generated)



[/sklearn.linear\\_model.LinearRegression.html](/sklearn.linear_model.LinearRegression.html))

- More on [model\\_evaluation](https://scikit-learn.org/stable/modules/model_evaluation.html) ([https://scikit-learn.org/stable/modules/model\\_evaluation.html](https://scikit-learn.org/stable/modules/model_evaluation.html))

Explain what's the meaning of reg.score return value in this notebook.

**A11** The meaning of reg.score in the notebook gives us the return coefficient of determination of the prediction. This highest value score produces is 1.0. This score comes from the residual sum over the total sum of squares which is then subtracted from 1. The closer to 1.0 the score is the better

```
In [109]: # A custom function to calculate r2 score
# Details on the custom scorers: https://scikit-learn.org/stable/modul

def r2score(ytrue, ypred):
    rss = ((ytrue - ypred)**2).sum()
    tss = ((ytrue - ytrue.mean()) ** 2).sum()
    r2 = 1 - rss/tss
    return r2

# Now do prediction on xtrain and xtest and check your r2 score by pri
trainpredict = reg.predict(xtrain)
testpredict = reg.predict(xtest)

print(r2score(ytrain, trainpredict), r2score(ytest, testpredict))

# You can see that reg.score values and your custom function output ar

0.8154769852171251 0.8449374024375086
```

One way of achieving linear regression is by minimizing the error between actual y and predicted y. The method is known as least square method. We will make our custom least square optimize to calculate model parameters that minimizes output error.

**Q12** Write a function which takes weights(or params), x and y and do following

- 1. calculate dot product between x and params , which is ypredicted
- 2. calculate difference between actual y and ypredicted
- 3. return the difference

**A12** complete the code below

```
In [113]: import scipy.optimize as optimization

def constraint(params, x, y):
    ypred = x@params
    return y-ypred

# Our initial params is a vector of size equal to dimension of x, or y
# You can create zeros vector using np.zeros(size)

# complete code
params = np.zeros(len(df.columns))
```

```
# Now study the documentation and complete following code
params, _ = optimization.leastsq(???, ????, ????)

# Now we have parameter or weight we can now create our model
model = lambda x:np.dot(x,params)

# Now predict ytrain using model and see first 5 predicted and actual
ypred_train = model(????)
# see first 5 predicted values
print(????)
# see first 5 actual values
print(????)

# Now predict ytest using model and see first 5 predicted and actual v
ypred_test = model(????)
print(????)
print(????)

# Now use custom made r2score calculator to calculate r2 score on both
print(r2score(ytest, ypred_test), r2score(ytrain, ypred_train))

# Add code to compare these results with the results you got from sklearn
??????
```

...