

VIETNAM NATIONAL UNIVERSITY, HO CHI MINH CITY  
UNIVERSITY OF TECHNOLOGY  
FACULTY OF COMPUTER SCIENCE AND ENGINEERING



## PROGRAMMING INTERGRATION PROJECT (CO3101) - CC01

---

Report

# Lasso Regression

---

Advisor: Assoc. Prof. Quản Thành Thơ

Students: Trịnh Mạnh Hùng - 1952740  
Nguyễn Quốc Việt - 1953096  
Trần Đức Nam - 1952861  
Nguyễn Quốc Văn Chương - 1950004  
Phạm Tấn Phuộc - 1952406

HO CHI MINH CITY, OCTOBER 2021



## Contents

<b>1 Overview Of The Lasso Regression</b>	<b>4</b>
1.1 Introduction . . . . .	4
1.2 "Best-Fit" In Lasso Regression . . . . .	6
1.2.1 MSE Fit Only . . . . .	7
1.2.2 MSE And L <sub>1</sub> Penalty . . . . .	7
1.3 Mathematical analysis of Linear regression . . . . .	9
1.3.1 Input and Output . . . . .	9
1.3.2 Loss function and Algebra solving . . . . .	9
1.4 Mathematical equation of Lasso Regression . . . . .	11
1.4.1 What is Regularization . . . . .	11
1.4.2 A large number of variables . . . . .	11
1.4.3 L <sub>1</sub> Regularization . . . . .	12
1.4.4 Cost function . . . . .	12
1.5 Algorithm . . . . .	13
1.5.1 Coordinatewise gradient descent (CGD) algorithm . . . . .	13
1.5.2 Gradient LASSO algorithm . . . . .	13
1.6 Advantages and Disadvantages of using Lasso Regression . . . . .	15
<b>2 Running Example</b>	<b>16</b>
2.1 Formulas and Given data . . . . .	16
2.2 Perform algorithm manually . . . . .	18
<b>3 Lasso Regression support tool</b>	<b>21</b>
<b>4 Metrics for evaluation</b>	<b>22</b>
4.1 Mean squared error (MSE)/ Root mean squared error (RMSE) . . . . .	22
4.2 Mean absolute error (MAE) . . . . .	22
4.3 R-squared . . . . .	22
<b>5 Toy problem demonstration</b>	<b>23</b>
5.1 Toy Dataset overview . . . . .	23
5.2 Running with Python code . . . . .	23
5.3 Implementation on web . . . . .	29
<b>6 Canada CO<sub>2</sub> Emisssions By Cars - Nam Trần</b>	<b>30</b>
6.1 Data Overview . . . . .	30
6.2 Implementation . . . . .	31
6.2.1 Data Preprocessing . . . . .	31
6.2.1.a Import Library . . . . .	31



6.2.1.b	Import Data . . . . .	31
6.2.1.c	Clean Missing Data . . . . .	32
6.2.1.d	Encoding Categorical . . . . .	33
6.2.1.e	Data Split . . . . .	36
6.2.2	Build Lasso Model . . . . .	36
6.2.3	Evaluation . . . . .	37
6.2.3.a	Lasso Evaluation . . . . .	37
6.2.4	Models Comparasion . . . . .	37
<b>7</b>	<b>COVID-19 vaccination vs mortality - Quốc Việt</b>	<b>39</b>
7.1	Data Overview . . . . .	39
7.2	Running with Python code . . . . .	40
7.3	Fitting model - Covid-19 Italy . . . . .	45
7.4	Evaluation . . . . .	46
7.5	Conclusion . . . . .	48
<b>8</b>	<b>Factors influencing Life Expectancy (WHO) - Tấn Phước</b>	<b>49</b>
8.1	Inspiration and Data Overview . . . . .	49
8.2	Data analysis . . . . .	49
8.2.1	Import library and data . . . . .	50
8.2.2	Data exploration . . . . .	51
8.2.3	Data preprocessing . . . . .	53
8.3	Fitting model . . . . .	55
8.3.1	Data splitting . . . . .	55
8.3.2	Train model . . . . .	55
8.3.3	Evaluation . . . . .	56
8.3.4	Comparison . . . . .	57
<b>9</b>	<b>Medical Cost Personal - Mạnh Hùng</b>	<b>59</b>
9.1	Background and Data Overview . . . . .	59
9.2	Implementation . . . . .	59
9.2.1	Import package and acquire the data . . . . .	59
9.2.2	Preprocessing data . . . . .	60
9.2.2.a	Check Null values . . . . .	60
9.2.2.b	Visualization . . . . .	61
9.2.2.c	Select features (based on correlation) . . . . .	62
9.3	Building the model and evaluation . . . . .	63
9.3.1	Data splitting . . . . .	63
9.3.2	Train model . . . . .	64
9.4	Evaluation and Comparision . . . . .	65



9.4.1	Evaluation . . . . .	65
9.4.2	Comparision . . . . .	65
<b>10</b>	<b>Predict final grade - Chương</b>	<b>68</b>
10.1	Data Overview . . . . .	68
10.2	Import libraries and preprocessing data . . . . .	68
10.2.1	Set up libraries and load data . . . . .	68
10.2.2	Preprocessing . . . . .	69
10.3	Visualization and selecting attributes . . . . .	71
10.3.1	Visualization . . . . .	71
10.3.2	Selecting attributes and split train/test . . . . .	73
10.4	Build and evaluate models . . . . .	74
10.4.1	Lasso Regression . . . . .	74
10.4.2	Random Forest and XGB . . . . .	75
10.5	Comparison and new prediction . . . . .	76
10.5.1	Compare 3 models . . . . .	76
10.5.2	New prediction value . . . . .	78
<b>11</b>	<b>Application</b>	<b>79</b>
11.1	Usage . . . . .	79
11.2	Evaluation . . . . .	79

## Abstract

Regression is one of popular task in machine learning that consists of mathematical methods that allow data scientists to predict a continuous outcome ( $y$ ) based on the value of one or more predictor variables ( $x$ ). This report discusses a traditional regression method called Lasso Regression algorithm, which is an advanced version of Linear Regression algorithm, from mathematical perspective. Additional, To make it easy, we'll use a lot of two dimensional space examples to illustrate the algorithm.

# 1 Overview Of The Lasso Regression

## 1.1 Introduction

Regression analysis is mainly used for prediction, forecasting, time series modeling, and determining the causal-effect relationship between variables.

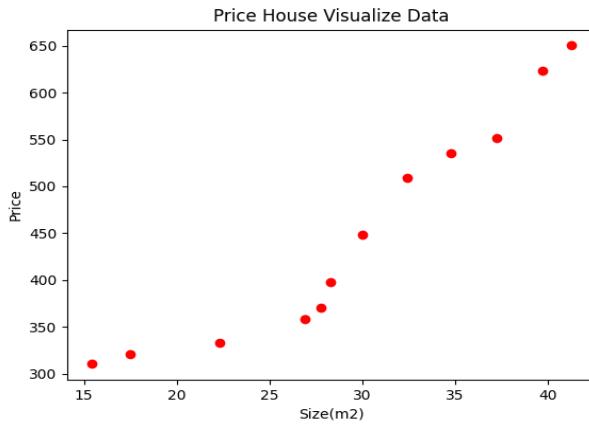
Lasso Regression or ‘Least Absolute Shrinkage and Selection Operator’, which is a supervised learning algorithms, do exactly what a Regression often do that it will try to plot a graph between the variables which “*best fits*” the given data points, using this plot, the machine learning model can make predictions about the data.

For examples, data points could be the collection of house price and their feature is size. Now, suppose that we have fifteen data points (or information of fifteen houses) as the following table:

Size( $m^2$ )	Price(USD)
15.4	310.5
17.5	320.7
22.3	333
26.9	358.5
27.8	370.3
28.3	398.2
30	448.5
32.4	509.2
34.8	535.1
37.2	551.4
39.7	623.4
41.3	650.2

Figure 1. Price House In England (2017)

Now, we plot this table into a two dimensional space with  $x$  axis represent the feature *size* and  $y$  axis represent *price* to have a clear look about the distribution of data points.



The goal of Lasso Regression now is to sketch a approximate line (with formula:  $h(x) = Ax + B$ ) to cross as close as possible all the points in the above two dimensional space.

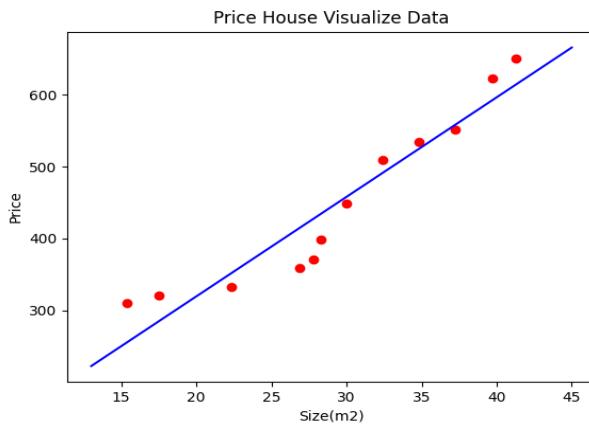


Figure 2. Lasso Regression Approximation Line

The blue line on the above image has the formula  $h(x) = 13.86x + 42.34$ . Basing on this blue line, whenever we has information about the size of a new house, we can predict the corresponding price of this house. For instance, if new house has the size is  $66\text{ m}^2$  ( $x = 66$ ), the price of this house should nearly equal 957 (USD).



## 1.2 "Best-Fit" In Lasso Regression

In the previous example, we see that the blue line is best approximation line to our data points. But, how can we know that this line is best fit?

One of the fundamentally basic idea is that using the MSE (Mean Square Error):

$$MSE = \frac{1}{m} \sum_{i=1}^m (y^{(i)} - h(x^{(i)}))^2$$

Here,

- $m$  is the total number of data points in the dataset.
- $y_{(i)}$  represents the value of target variable for  $i^{th}$  training example.
- $h(x_{(i)})$  represents the value basing on the formula of prediction line with corresponding  $x^{(i)}$

In the other words, instead of drawing many lines and check to see that whether this line is satisfy our requirement in the two dimensional space or not , we now find the value of **Weight-A**, and **bias-B** (in function  $h(x) = Ax + B$ ) to minimize the above function MSE. However, in Lasso regression, it is said that not only use the MSE as way to find the blue line, we also add one more constraint, called  $L_1$  penalty to make our model be more regularize

$$L_1 = \frac{1}{m} \lambda \sum_{(j=1)}^n w_j$$

where

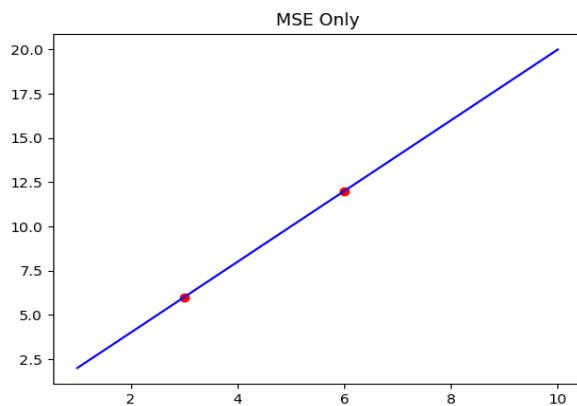
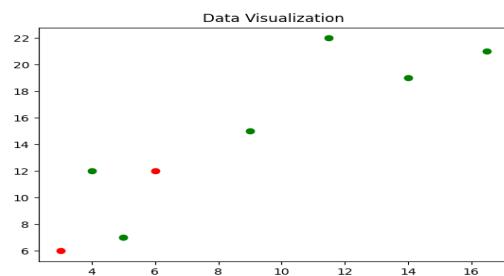
- $\lambda$  is the regularization strength.
- $w_{(j)}$  represents the weight for  $j^{th}$  feature.
- $n$  is the number of features in the dataset.
- $m$  is the total number of data points in the dataset.

**Finally**, we have function and our goal is to minimize it with respect to weight and bias. In additionally, we will call this function is **Objective Function** or **Loss Function** through the report from now.

$$Loss = \frac{1}{m} \left[ \sum_{i=1}^m (y^{(i)} - h(x^{(i)}))^2 + \lambda \sum_{(j=1)}^n |w_j| \right]$$

To illustrate how the  $L_1$  penalty works, we will consider bellow example.

Suppose that the red points are the points we use to sketch the approximation line and the green ones are to test this line. Next, we separate into two scenarios. First scenario, we apply only the MSE (**Linear Regression**) function to find the approximation line. The other one, we apply both MSE and  $L_1$  penalty (**Lasso Regression**), to see the differences.



### 1.2.1 MSE Fit Only

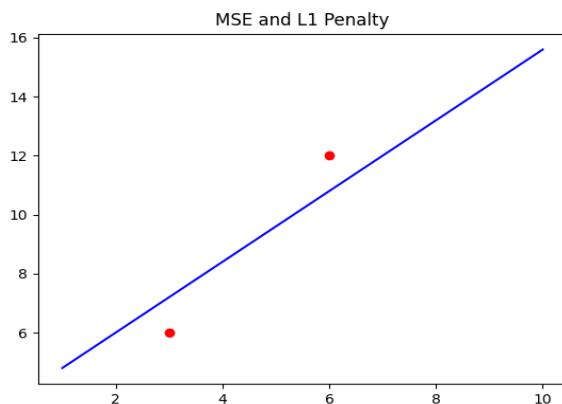
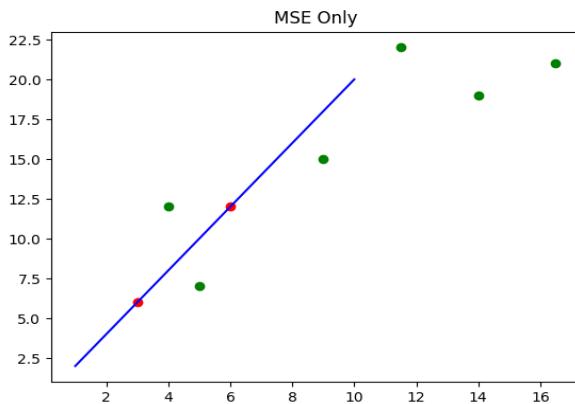
If we use only the MSE without any constraints, the approximation line will be completely fitted to our two training points ( $\text{MSE} = 0$ ).

However, when we test this blue line by our testing set (green data points), the line is not good as well as we saw anymore. In the other words, the line creates a large of variance.

In machine learning lingo, we'd say that the above blue line is **Over Fit** to the Training Data (red points).

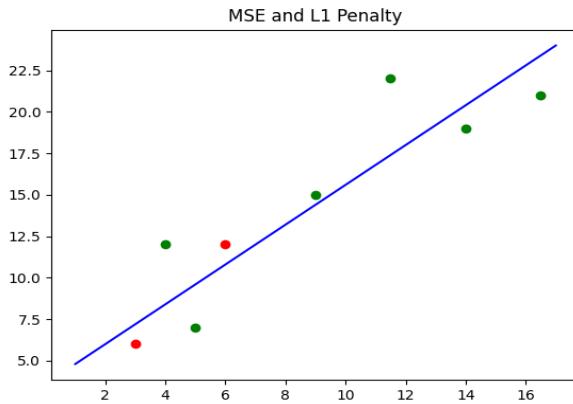
### 1.2.2 MSE And $L_1$ Penalty

Now, instead of using only MSE method, we use simultaneously  $L_1$  penalty. Then the blue line will become



Obviously, the second scenario approximation line has the smaller slope than the first scenario . In the other words, the second line does not fit the Training Data as well but in return for that changes, we get a significant drop in **Variance**.

So, by starting with a slight worse fit,  $L_1$  penalty can provide better long terms predictions. Conclusion, the main idea behind the  $L_1$  penalty is to make the slope of the approximation line be steep, then the prediction for output  $y$  is very sensitive to relative small changes in  $x$ .



## 1.3 Mathematical analysis of Linear regression

### 1.3.1 Input and Output

The Linear Regression Algorithms take a set of Observations  $X = \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(n)} \end{bmatrix} \in R^{nx(m+1)}$  where each observation is a  $m$ -dimensional row vector,  $n$  is the number of observations (members) and  $m$  is the numbers of features in each observation. **Note:** To apply the algorithm,  $n > m + 1$ .

The algorithm outputs is a prediction vector  $\hat{y} = \begin{bmatrix} \hat{y}^{(1)} \\ \hat{y}^{(2)} \\ \vdots \\ \hat{y}^{(n)} \end{bmatrix} \in R^{n*1}$  which is the approximate result of the expected output.

### 1.3.2 Loss function and Algebra solving

Before we start, let's understand some important notation.

$w = [w_0, w_1, w_2, \dots, w_m]^T$  is a coefficient vector that needs to be optimized,  $w_0$  is often called as bias.

$x^{(i)} = [1, x_1^{(i)}, x_2^{(i)}, \dots, x_m^{(i)}]$  is the observation  $i^{th}$  in the set of  $n$  observations, each observation has  $m$  features.

Because our model is linear relation  $\Rightarrow \hat{y}_i = x^{(i)}w$

General form of Linear Regression model:

$$X \in R^{n*(m+1)}, w \in R^{(m+1)*1}, y \in R^{n*1}$$



$$X = \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(n)} \end{bmatrix}, w = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_m \end{bmatrix}, \hat{y} = Xw = \begin{bmatrix} w_0 + w_1x_1^{(1)} + \dots + w_mx_m^{(1)} \\ w_0 + w_1x_1^{(2)} + \dots + w_mx_m^{(2)} \\ w_0 + w_1x_1^{(3)} + \dots + w_mx_m^{(3)} \\ \vdots \\ w_0 + w_1x_1^{(n)} + \dots + w_mx_m^{(n)} \end{bmatrix}, y - \hat{y} = \begin{bmatrix} y^{(1)} - (w_0 + w_1x_1^{(1)} + \dots + w_mx_m^{(1)}) \\ y^{(2)} - (w_0 + w_1x_1^{(2)} + \dots + w_mx_m^{(2)}) \\ y^{(3)} - (w_0 + w_1x_1^{(3)} + \dots + w_mx_m^{(3)}) \\ \vdots \\ y^{(n)} - (w_0 + w_1x_1^{(n)} + \dots + w_mx_m^{(n)}) \end{bmatrix}$$

Loss function:  $L = \frac{1}{2} * \frac{1}{n} * \sum_{i=1}^n (y^i - x^i w)^2$

Definition of Euclidean norm or norm2:  $\|z\|_2 = (z_1^2 + z_2^2 + \dots + z_n^2)^{\frac{1}{2}} \implies \|z\|_2^2 = (z_1^2 + z_2^2 + \dots + z_n^2)$

$$\implies L = \frac{1}{2} * \frac{1}{n} \|y - \hat{y}\|_2^2 = \frac{1}{2} * \frac{1}{n} * (y - \hat{y})^T * (y - \hat{y})$$

Now, we have already know the general formula of the loss function  $L$  Linear Regression in term of Linear Algebra. To completely solving the Linear Regression problems, we must try to find the global minimum of loss function  $L$  in which the parameters  $w$  is unknown. In the other words, we must find the optimized parameters  $w$  which is a vector such that the value of  $L$  is minimum.

However, the min value of Loss function  $L$  will be the similar with  $\frac{L}{N}$  so we'll find the minimum of function  $L_1$ :

$$L_1 = \frac{1}{2} * (y - \hat{y})^T * (y - \hat{y})$$

Take the derivative on both sides of function  $L_1$  with respect to  $w$  we have:

$$\begin{aligned} &\Rightarrow \frac{d(L)}{d(w)} = \frac{d(\frac{1}{2} * (y - \hat{y})^T * (y - \hat{y}))}{d(w)} \\ &\Leftrightarrow \frac{d(L_1)}{d(w)} = \frac{d(\frac{1}{2} * (y - (Xw))^T * (y - Xw))}{d(w)} \\ &\Leftrightarrow \frac{d(L_1)}{d(w)} = \frac{1}{2} * 2 * (-X^T) * (y - Xw) \\ &\Leftrightarrow \frac{d(L_1)}{d(w)} = X^T X w - X^T y \end{aligned}$$

To find the minimum, let the left hand side be zero:

$$\begin{aligned} &\Leftrightarrow 0 = X^T X w - X^T y \\ &\Leftrightarrow X^T X w = X^T y \end{aligned}$$

Assume that,  $X^T X$  is invertible matrix:

$$\begin{aligned} &\Leftrightarrow (X^T X)^{-1} (X^T X) w = (X^T X)^{-1} (X^T y) \\ &\Leftrightarrow I w = (X^T X)^{-1} (X^T y) \\ &\Leftrightarrow w = (X^T X)^{-1} (X^T y) \end{aligned}$$

Conclusion, the Loss function  $L_1$  or  $L$  reach the minimum at  $w = (X^T X)^{-1} (X^T y)$  if and only if  $X^T X$  is invertible.

## 1.4 Mathematical equation of Lasso Regression

### 1.4.1 What is Regularization

Regularization solves the problem of overfitting. Overfitting causes low model accuracy. It happens when the model learns the data as well as the noises in the training set.

Noises are random datum in the training set which don't represent the actual properties of the data.

$$Y \approx C_0 + C_1 X_1 + C_2 X_2 + \dots + C_p X_p$$

Y represents the dependent variable, X represents the independent variables and C represents the coefficient estimates for different variables in the above linear regression equation.

A loss function called the sum of squares is used to fit the model. The coefficients in the equation are chosen so that the loss function is reduced to the smallest possible value. If there is a lot of irrelevant data in the training set, the wrong coefficients will be chosen. This will result in bad model predictions.

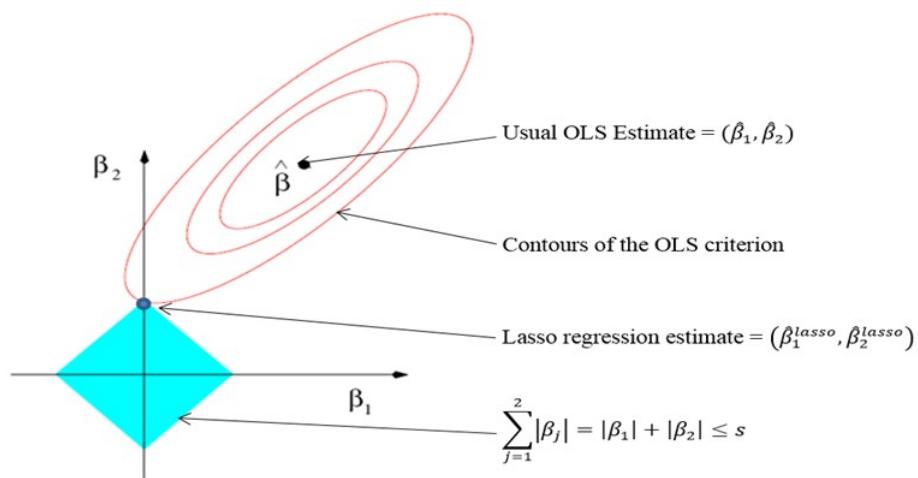
In cases like this, we can use regularization to regularize or shrink these wrongly learned coefficients to zero. Lasso regression is one of the popular techniques used to improve model performance.

### 1.4.2 A large number of variables

A large number of variables means:

1. The large number here means that the model tends to over-fit. Theoretically, a minimum of ten variables can cause an overfitting problem.
2. When you face computational challenges due to the presence of n number of variables. Although, given today's processing power of systems, this situation arises rarely.

The following diagram is the visual interpretation comparing OLS and lasso regression:





The LASSO is not very good at handling variables that show a correlation between them and thus can sometimes show very wild behavior.

#### 1.4.3 L1 Regularization

Lasso regression performs L1 regularization, which adds a penalty equal to the absolute value of the magnitude of coefficients. This type of regularization can result in sparse models with few coefficients

Some coefficients can become zero and eliminated from the model. Larger penalties result in coefficient values closer to zero, which is the ideal for producing simpler models.

Our regularized model may have a slightly high bias than linear regression but less variance for future predictions.

Residual Sum of Squares +  $\lambda$  \* (Sum of the absolute value of the magnitude of coefficients)

$$\sum_{i=1}^n (y_i - \sum_j x_{ij}\beta_j)^2 + \lambda \sum_{j=1}^p |\beta_j|$$

Where:

- $\lambda$  denotes the amount of shrinkage, which controls the strength of the L1 penalty
- $\lambda = 0$  implies all features are considered and it is equivalent to the linear regression where only the residual sum of squares is considered to build a predictive model
- $\lambda = \infty$  implies no feature is considered i.e, as  $\lambda$  closes to infinity it eliminates more and more features
- The bias increases with increase in  $\lambda$
- Variance increases with decrease in  $\lambda$

In lasso, the coefficients which are responsible for large variance are converted to zero.

#### 1.4.4 Cost function

Lasso Regression or ('Least Absolute Shrinkage and Selection Operator') also works with an alternate cost function:

$$J(W) = \frac{1}{2N} \sum_{i=1}^N ((W_0 + W_1 X_1^{(i)} + \dots + W_P X_P^{(i)}) - Y_i)^2 + \frac{\lambda}{2N} \sum_{j=1}^P |W_j|$$

The derivative of this cost function has no closed form (due to the L1 loss on the weights), which means we can not simply apply gradient descent. Lasso allows for the possibility that a coefficient can actually be forced to zero, essentially making Lasso a method of model selection as well as a regression technique.



## 1.5 Algorithm

LASSO is computationally demanding since the L1 constraint is not differentiable. In this section, we will talk about some algorithms related to LASSO Regression.

### 1.5.1 Coordinatewise gradient descent (CGD) algorithm

Recall that for LASSO we need to solve:  $\text{Minimize}_{\beta} R(\beta)$  subject to  $\|\beta\|_1 \leq \lambda$ .  
For simplicity, we let  $\lambda = 1$ . For other  $\lambda$ , we rescale the inputs by multiplying  $\lambda$ .

Let  $e_k$  be the p-dimensional vector whose k-th entry is 1 and the others are zero (For example:  $e_1 = (1, 0, \dots, 0)$ ) (called it an “coordinate vector”).

Let  $\varepsilon = e_k, -e_k : k = 1, \dots, p$ .

The LASSO problem can be restated as  $\text{Minimize}_{\beta(\varepsilon)} R(\beta)$  Where  $\text{co}(\varepsilon)$  is the convex hull of  $\varepsilon$ .

The CGD algorithm updates  $\beta$  as follows.

- Get the gradient vector  $R^{(1)}(\beta) = \frac{\partial R}{\partial \beta}$ .
- Choose the coordinate vector  $e \in \varepsilon$  which minimizes  $e'R^{(1)}(\beta)$ .
- Find the convex combination of  $\beta$  and  $e$  which minimizes  $R$ , which means finding  $\alpha \in [0, 1]$  that minimizes  $R(\alpha\beta + (1 - \alpha)e)$ .
- Update  $\beta = \alpha\beta + (1 - \alpha)e$ .

In brief:

- The LASSO problem is an optimization problem over the simplex of a given set of vertices.
- The CGD algorithm repeatedly finds the optimal vertex (i.e smallest gradient) and takes the optimal convex combination (i.e. minimizing  $R$ ).

### 1.5.2 Gradient LASSO algorithm

Let  $A = \emptyset$  ( $A$ : active set).

The gradient LASSO algorithm consists of the two steps: addition and deletion.

- Addition: Move the solution using the CGD algorithm, which may result in adding a new coordinate vector into  $A$ ..
- Deletion: Move toward the (projected) gradient direction on the simplex of  $A$  (i.e.  $\text{co}(A)$ ), which may result in deletion of a coordinate vector from  $A$ .

The gradient LASSO algorithm repeats the addition and deletion steps until the solution converges.

We should note that the gradient LASSO algorithm always converges faster than the CGD algorithm since the deletion step decreases the empirical risk.

**Remark:**



- The gradient LASSO algorithm does not require matrix inversion.
- In each iteration, we need two one-dimensional optimizations - one for addition and the other for deletion, which can be done easily.
- Hence, the gradient LASSO algorithm has all of the advantages of the CGD algorithm, and at the same time, it improves the speed of convergence significantly.



## 1.6 Advantages and Disadvantages of using Lasso Regression

By using Lasso Regression to solve some specific problems in the real world, we will need to first consider its advantages. LASSO isn't a type of regression, it is a regularization or a method of model building and variable selection that can be applied to many types of regression, including ordinary least square methods, logistic regression, and . . .

The biggest pro of LASSO is that it is better than the original methods of automatic variable selection such as forward, backward and stepwise. These listing methods may result in giving wrong result, the results got from LASSO are better in general.

### LASSO gives rise to the following merits:

- Playing a roll of regularization method, LASSO can avoid overfitting. It can be applied even when number of features is in a higher number or larger than the number of data.
- It can also do feature selection.<sup>1</sup>
- LASSO is fast in terms of inference and fitting the model.

Looking at the disadvantages, LASSO problems is the idea to be automatic. When LASSO performs a regularization on the data, it avoid involving thinking or in another word, it can produce models that clearly make no sense.

### LASSO regression may show the demerits:

- The model selected by LASSO is not always stable. Sometimes, on different bootstrapped data, the feature selection performed can be very different.
- LASSO, as a method for regression, can ignores the variables that is not significant, some of these variables may, however, be interesting or play an important roll.
- When it comes to highly correlated features, LASSO may randomly select one of them or a part of them.

---

<sup>1</sup>Feature selection by LASSO

## 2 Running Example

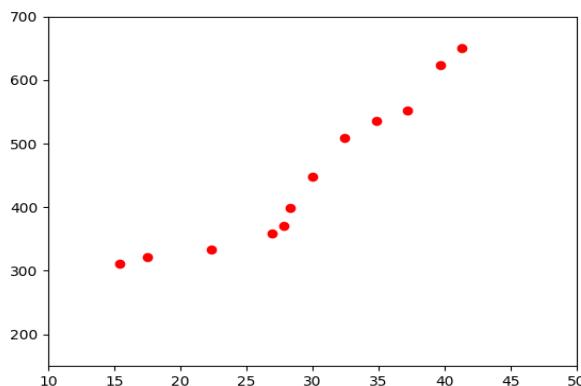
### 2.1 Formulas and Given data

To illustrate Lasso regression, we are going to use 12 data points in this table of data, Figure 3:

Size( $m^2$ )	Price(USD)
15.4	310.5
17.5	320.7
22.3	333
26.9	358.5
27.8	370.3
28.3	398.2
30	448.5
32.4	509.2
34.8	535.1
37.2	551.4
39.7	623.4
41.3	650.2

Figure 3. Price House In England (2017)

For better visualizing, we will plot these data points on a 2D space:

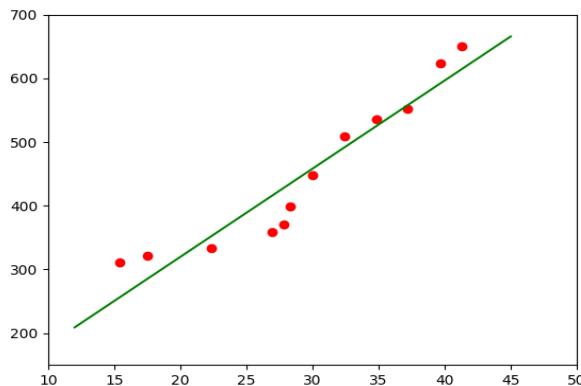


What we need to do here is finding a line that can cross these data points with a minimum Mean Squared Error. To do this task, we will use two ways for performing Lasso Regression on this dataset. The first method is using scikit-learn, which is a Machine Learning library, to train a Lasso Regression model. Following that, we will use gradient descent as the second method to fit our Lasso Regression model.

Let's begin with our first method. With scikit-learn, we just need few lines of code to train a Lasso Regression model, which is as below:

```
1 import numpy as np
2 import matplotlib
3 import matplotlib.pyplot as plt
4 from sklearn import linear_model
5 # Data
6 X = np.array([[15.4, 17.5, 22.3, 26.9, 27.8, 28.3, 30, 32.4, 34.8, 37.2, 39.7,
7 41.3]]).T
7 Y = np.array([[310.5, 320.7, 333, 358.5, 370.3, 398.2, 448.5, 509.2, 535.1, 551.4,
623.4, 650.2]]).T
8 # Draw data
9 fig1 = plt.figure("Lasso Regression")
10 ax = plt.axes(xlim=(10,50), ylim=(150,700))
11 plt.plot(X, Y, 'ro')
12
13 lr = linear_model.Lasso(alpha=1.6, max_iter=100, tol=0.1, )
14 lr.fit(X,Y)
15 x0_gd = np.linspace(12,45,2)
16 y0_sklearn = lr.intercept_[0] + lr.coef_[0]*x0_gd
17
18 plt.plot(x0_gd,y0_sklearn, color="green")
19
20 plt.show()
```

The output:



With just few lines of code, we now have a line (green line) fit very well with our dataset. Scikit-learn is surely convenient.

However, to understand how Lasso Regression really works, we are going to work manually with this model through Gradient Descent algorithm. Basically, Gradient Descent is an optimization algorithm for finding a local minimum of a differentiable function. It is simply used in machine learning to find the values of a function's parameters (coefficients) that minimize a cost function as far as possible. The idea is to take repeated steps in the opposite direction of the gradient (or approximate gradient) of the function at the current point, because this is the direction of steepest descent. Based on this algorithm, we are able to fit our Lasso Regression model after some iterations. Now, let's get started.



Initially, we have the loss function:

$$\text{Loss} = \frac{1}{m} \left[ \sum_{i=1}^m (y^{(i)} - h(x^{(i)}))^2 + \lambda \sum_{j=1}^n |w_j| \right] = \frac{1}{m} \left[ \|Xw - y\|_2^2 + \lambda \|w\|_1 \right]$$

Here,  $X \in R^{m*n}$ ,  $w \in R^{n*1}$ ,  $y \in R^{m*1}$ ,  $m$  is the number of observations,  $n$  is number of features in each observation.

According to the idea of Gradient Descent, we need to differentiate this loss function. However, as we mentioned before

$$\text{Loss} = \text{MSE} + L_1 \text{ Penalty}$$

So, differentiating Loss function we have:

$$\begin{aligned} \frac{d(\text{Loss})}{d(w)} &= \frac{1}{m} \times \frac{d(\|Xw - y\|_2^2)}{d(w)} + \frac{\lambda}{m} \times \frac{d(\|w\|_1)}{d(w)} \\ &= \frac{1}{m} \times X^T(Xw - y) + \frac{\lambda}{m} \times \begin{bmatrix} \frac{d(|w_1|+|w_2|+\dots+|w_n|)}{d(w_1)} \\ \frac{d(|w_1|+|w_2|+\dots+|w_n|)}{d(w_2)} \\ \vdots \\ \frac{d(|w_1|+|w_2|+\dots+|w_n|)}{d(w_n)} \end{bmatrix} \end{aligned} \quad (1)$$

Using the formula that we have found above, we can find the derivative of Loss function at a certain point. With the value of derivative, we are able to determine the right way to update our coefficients - w.

## 2.2 Perform algorithm manually

Gradient descent is an algorithm to find the minimum of Loss function based on differentiation. The algorithm is as follows:

1. Initialize value for vector  $u_0 = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ \vdots \\ w_n \end{bmatrix}$  randomly. **Note:** the vector  $w \in R^{n*1}$  will be renamed to vector u.
2. Assign  $u_i = u_{(i-1)} - \text{learning\_rate} * \frac{d(\text{Loss})}{d(u)}$  (learning\_rate is a positive constant, for instance: learning\_rate = 0.001)
3. Calculating  $\text{Loss}(u_i)$  again: If  $\text{Loss}(u_i)$  is small enough then we will terminate our algorithm; otherwise, we repeat with step 2.



According to Figure 3, we also have  $X = \begin{bmatrix} 1.0 & 15.4 \\ 1.0 & 17.5 \\ 1.0 & 22.3 \\ 1.0 & 26.9 \\ 1.0 & 27.8 \\ 1.0 & 28.3 \\ 1.0 & 30.0 \\ 1.0 & 32.4 \\ 1.0 & 34.8 \\ 1.0 & 37.2 \\ 1.0 & 39.7 \\ 1.0 & 41.3 \end{bmatrix}$  and  $y = \begin{bmatrix} 310.5 \\ 320.7 \\ 333.0 \\ 358.5 \\ 370.3 \\ 398.2 \\ 448.5 \\ 509.2 \\ 535.1 \\ 551.4 \\ 623.4 \\ 650.2 \end{bmatrix}$ . Also, since our dataset have 12 datapoints, so  $m = 12$ . Assume that  $\lambda = 1.8$  and  $learning\_rate = 0.001$

Now, Let's solve the above example by hand.

**Step 1:** Initialize value for vector  $u_0 = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$ .

**Step 2:**  $u_1 = u_0 - 0.001 * \frac{d(Loss)}{d(u)} = u_0 - 0.001 * \left( \frac{1}{m} * X^T(Xu_0 - y) + \frac{\lambda}{m} * \begin{bmatrix} \frac{d(|w_1|+|w_2|)}{d(w_1)} \\ \frac{d(|w_1|+|w_2|)}{d(w_2)} \end{bmatrix} \right)$ .

Substituting X, y, m , n,  $u_0$ . we have:  $u_1 = \begin{bmatrix} 2.04 \\ 2.31 \end{bmatrix}$

**Step 3:** Calculating  $Loss(u_1) = \frac{1}{m} \left[ \|Xu - y\|_2^2 + \lambda \|u\|_1 \right] = 77155.83$ . This result is too large. we want to the value of Loss should be less than 700. So, we'll continue update vector  $u_1$  to reduce the value of Loss.

**Step 4:**  $u_2 = u_1 - 0.001 * \frac{d(Loss)}{d(u)} = u_1 - 0.001 * \left( \frac{1}{m} * X^T(Xu_1 - y) + \frac{\lambda}{m} * \begin{bmatrix} \frac{d(|w_1|+|w_2|)}{d(w_1)} \\ \frac{d(|w_1|+|w_2|)}{d(w_2)} \end{bmatrix} \right)$ .

Substituting X, y, m , n,  $u_1$ . we have:  $u_2 = \begin{bmatrix} 2.079 \\ 3.508 \end{bmatrix}$

**Step 5:**  $Loss(u_2) = 63572.691 > 700$ . So, continuing update  $u_2$  to reduce the value of Loss.

**Step 6:**  $u_3 = u_2 - 0.001 * \frac{d(Loss)}{d(u)} = \begin{bmatrix} 2.114 \\ 4.590 \end{bmatrix}$

**Step 7:**  $Loss(u_3) = 52402.373 > 700$ . So, continuing update  $u_3$  to reduce the value of Loss.

$\vdots$

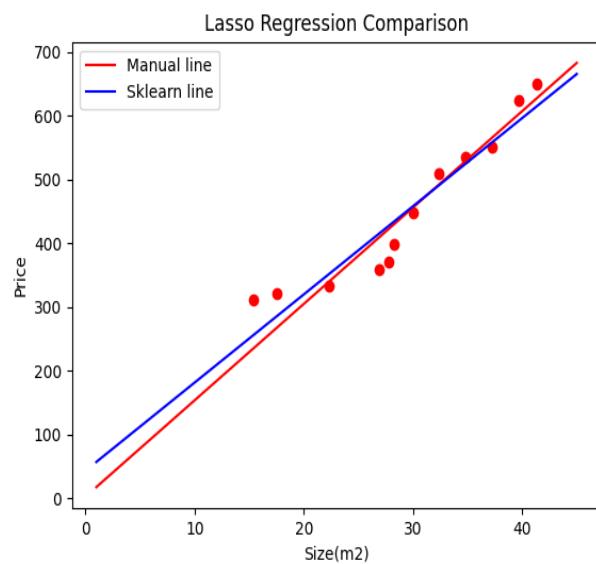
**Step 119:**  $Loss(u_{39}) = 705.743 > 700$ . So, continuing update  $u_{39}$  to reduce the value of Loss.

**Step 120:**  $u_{40} = u_{39} - 0.001 * \frac{d(Loss)}{d(u)} = \begin{bmatrix} 2.470 \\ 15.119 \end{bmatrix}$

**Step 121:**  $Loss(u_{40}) = 699.976 < 700$  (satisfied). That means our line will have formula:  $y_1 = 15.119x + 2.470$

**Conclusion,** Our manual addressing stop after exactly 121 steps. And when compared to the formula which was created by sklearn,  $y_2 = 13.929x + 44.193$  we can calculate the Loss value of between two line  $y_1$  and  $y_2$  are respectively 688.976 and 646.331 respectively. So, our handling result is acceptable.

Finally, let take a look at how the line  $y_1$  and  $y_2$  are plotted in the two dimensional space.





### 3 Lasso Regression support tool

In Python, Scikit Learn supports users to build Lasso regression models using `sklearn.linear_model.Lasso` library.

- class `sklearn.linear_model.Lasso(alpha=1.0, *, fit_intercept=True, normalize='deprecated', precompute=False, copy_X=True, max_iter=1000, tol=0.0001, warm_start=False, positive=False, random_state=None, selection='cyclic')`

Main Parameters:

- **alpha****float, default=1.0:** Constant that multiplies the L1 term. Defaults to 1.0. alpha = 0 is equivalent to an ordinary least square, solved by the LinearRegression object. For numerical reasons, using alpha = 0 with the Lasso object is not advised. Given this, you should use the LinearRegression object.
- **fit\_intercept****bool, default=True:** Whether to calculate the intercept for this model. If set to False, no intercept will be used in calculations (i.e. data is expected to be centered).
- **normalize****bool, default=False:** This parameter is ignored when `fit_intercept` is set to False. If True, the regressors X will be normalized before regression by subtracting the mean and dividing by the l2-norm. If you wish to standardize, please use StandardScaler before calling fit on an estimator with `normalize=False`.
- **precompute****'auto', bool or array-like of shape (n\_features, n\_features), precompute:** Whether to use a precomputed Gram matrix to speed up calculations. The Gram matrix can also be passed as argument. For sparse input this option is always False to preserve sparsity.
- **copy\_X****bool, default=True:** If True, X will be copied; else, it may be overwritten.
- **max\_iter****int, default=1000:** The maximum number of iterations.

Example:

```
1  >>> from sklearn import linear_model
2  >>> clf = linear_model.Lasso(alpha=0.1)
3  >>> clf.fit([[0,0], [1, 1], [2, 2]], [0, 1, 2])
4      Lasso(alpha=0.1)
5  >>> print(clf.coef_)
6      [0.85 0.   ]
7  >>> print(clf.intercept_)
8      0.15...
```

Methods:

- **fit(X, y, sample\_weight=None, check\_input=True):** Fit model with coordinate descent.
- **get\_params(deep=True):** Get parameters for this estimator.
- **predict(X):** Predict using the linear model.
- **score(X, y, sample\_weight=None):** Return the coefficient of determination of the prediction



## 4 Metrics for evaluation

There are three common ways to evaluate a regression models:

### 4.1 Mean squared error (MSE)/ Root mean squared error (RMSE)

Mean Square Error is an absolute measure of the goodness for the fit. MSE is calculated by the sum of square of prediction error which is real output minus predicted output and then divide by the number of data points. It gives you a real number to compare against other model results and help you select the best regression model.

$$MSE = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2 \quad (2)$$

Root Mean Square Error(RMSE) is the square root of MSE. It is used more commonly than MSE because firstly sometimes MSE value can be too big to compare easily.

### 4.2 Mean absolute error (MAE)

Mean Absolute Error(MAE) is similar to Mean Square Error(MSE). However, instead of the sum of square of error in MSE, MAE is taking the sum of the absolute value of error. MSE gives larger penalization to big prediction error by square it while MAE treats all errors the same.

$$MAE = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}_i| \quad (3)$$

### 4.3 R-squared

R-Square (square of Correlation Coefficient(R)) measures how much variability in dependent variable can be explained by the model. R Square is calculated by the sum of squared of prediction error divided by the total sum of the square which replaces the calculated prediction with mean. R Square value is between 0 to 1 and a bigger value indicates a better fit between prediction and actual value. Moreover, R-squared does not take into consideration of overfitting problem.

$$R^2 = 1 - \frac{SS_{\text{Regression}}}{SS_{\text{Total}}} = 1 - \frac{\sum_i (y_i - \hat{y}_i)^2}{\sum_i (y_i - \bar{y})^2} \quad (4)$$



## 5 Toy problem demonstration

### 5.1 Toy Dataset overview

In the future, let's say we will become a real estate agent, and we are going in charge of selling a new house. We don't know the price of the house we want to sell, and we want to infer it by comparing it with other houses. We will look at features of the house which could influence the house, such as size, number of rooms, location, crime rate, school quality, distance to commerce, etc. At the end of the day, what we want is that a formula on all these features which gives us the price of the house, or at least an estimate for it. That is the idea spreading through housing prices problem, for the introduction with the Regression type model, this is the reason give our group idea to choose this dataset.

We are going to take a overview about this dataset. Since house price is a continues variables, this is so perfect for a regression model in general or just only the idea for LASSO. The data contains the following columns:

- **Avg. Area Income:** The average Income of the residents of the city that the house is located in
- **House Age:** The average age of houses in the same city
- **Number of Rooms:** The average number of rooms of houses in the same city
- **Number of Bedrooms:** The average number of bedrooms of houses in the same city
- **Area Population:** The population of city that the house is located in
- **Price:** The price of the house

### 5.2 Running with Python code

Our group will present the toy demonstration with some basic steps namely as follow:

- Data cleaning
- Selection of attributes
- Plotting to have an overview at first
- Perform algorithm
- Evaluation (this part will be update later on follow the guidance of Professor)

Our group choose to use Google Colab as it is easier to perform teamworking and sharing knowledge. By Google Colab, the code will be separate into many Code sections that can run separately, the benefit can be listed simple here is that we can recompile just a small section of codes so that it can reduce the compile time compare to the fully code.



The first step is to import some library that involve in the demonstration.

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import pandas as pd
4 import seaborn as sns
5 from sklearn.linear_model import Lasso
```

The library **numpy** is a Python library used for working with arrays. It also has functions for working in domain of linear algebra, fourier transform, and matrices.

The part of **matplotlib**, matplotlib.pyplot is a collection of functions that make matplotlib work like MATLAB. Each pyplot function makes some change to a figure: e.g., creates a figure, creates a plotting area in a figure, plots some lines in a plotting area, decorates the plot with labels, etc. The library **pandas** is a fast, powerful, flexible and easy to use open source data analysis and manipulation tool, built on top of the Python programming language.

The library **seaborn** is a Python data visualization library based on matplotlib. It provides a high-level interface for drawing attractive and informative statistical graphics.

Finally, the tool for our model is imported from **sklearn** - a Python module for machine learning built on top of SciPy and is distributed under the 3-Clause BSD license, with the **linear\_model** of **Lasso**.

To import the data, we will follow the id of the link on Google Drive. This makes the dataset accessible by many people in order to view it, just need to follow the ID.

```
1 # ID data: 1tjyJB0DrRom0b5E5ZKoFI6azSEZIkzwi
2 %%shell
3 cd /content
4 gdown -q --id "1tjyJB0DrRom0b5E5ZKoFI6azSEZIkzwi"
```

Next, we check if the data is imported successfully, if Yes then print 5 rows in the head and 5 rows in the tail, check it with our own dataset to ensure that it is correct.

```
1 df = pd.read_csv("House_price.csv")
2 print(df.head())
3 print(df.tail())
```

## Output

```
[3] df = pd.read_csv("House_price.csv")
print(df.head())
print(df.tail())

      Avg. Area Income House Age ... Area Population      Price
0    79545.45857   5.682861   ...  23086.80050  1.059034e+06
1    79248.64245   6.002900   ...  40173.07217  1.505891e+06
2    61287.06718   5.865890   ...  36882.15940  1.058988e+06
3    63345.24005   7.188236   ...  34310.24283  1.260617e+06
4    59982.19723   5.040555   ...  26354.10947  6.309435e+05

[5 rows x 6 columns]
      Avg. Area Income House Age ... Area Population      Price
4543  84556.63627   5.092459   ...  31797.31744  1326846.699
4544  82732.98111   5.332712   ...  41136.52722  1252663.621
4545  65694.05127   6.436741   ...  43406.71203  1359762.694
4546  76061.35071   7.148713   ...  40876.96459  1778013.334
4547  66935.47508   6.376390   ...  17573.61710  1109059.054

[5 rows x 6 columns]
```



Next, we will get deeper view into the data by using the `df.info()`, this function will return all the columns of the dataset together with the Data type (Dtype).

### Output

```
[ ] df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4548 entries, 0 to 4547
Data columns (total 6 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   Avg. Area Income    4548 non-null   float64
 1   House Age          4548 non-null   float64
 2   Number of Rooms    4548 non-null   float64
 3   Number of Bedrooms 4548 non-null   float64
 4   Area Population    4548 non-null   float64
 5   Price              4548 non-null   float64
dtypes: float64(6)
memory usage: 213.3 KB
```

With the toy problem we have here, it is so perfect that all the rows in the Dataset is non-null and so perfect that the Data type is float64. This results in the benefits that we don't have to perform some changes to our data to best fit the model. The Data cleaning process is also covered by the non-null, we can also check the number of null value by using another function of **pandas: df.isnull()**.

The Selection of attributes step is performed after we get the description about the dataset with the supported function `df.describe()`:

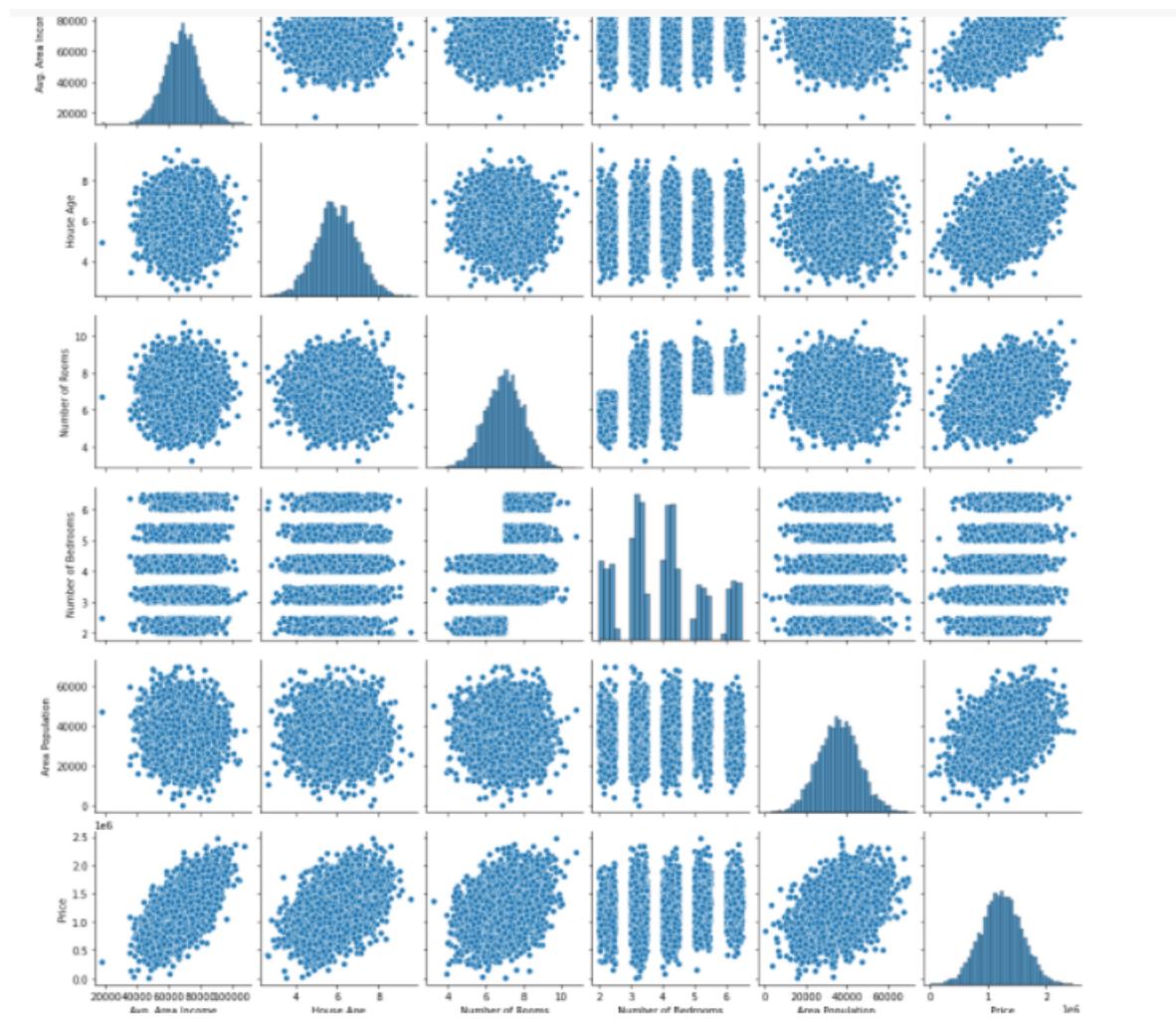
### Output

```
[6] df.describe()

      Avg. Area Income  House Age  Number of Rooms  Number of Bedrooms  Area Population        Price
count    4548.000000  4548.000000  4548.000000  4548.000000  4548.000000  4.548000e+03
mean    68611.700818   5.978918   6.987646   3.981693  36187.469334  1.233916e+06
std     10686.487761   0.990850   1.006587   1.230939  9910.189915  3.545676e+05
min     17796.631190   2.644304   3.236194   2.000000  172.610686  1.593866e+04
25%     61485.150193   5.332187   6.299692   3.140000  29423.163510  9.977751e+05
50%     68817.036575   5.960872   7.002245   4.050000  36215.560985  1.234571e+06
75%     75820.741747   6.658368   7.665871   4.490000  42880.554642  1.470616e+06
max     107701.748400  9.519088  10.759588   6.500000  69592.040240  2.469066e+06
```

Next, data visualization, to make it easy, our group chooses to plot the Pairplot by the command `sns.pairplot(df)`:

#### Output

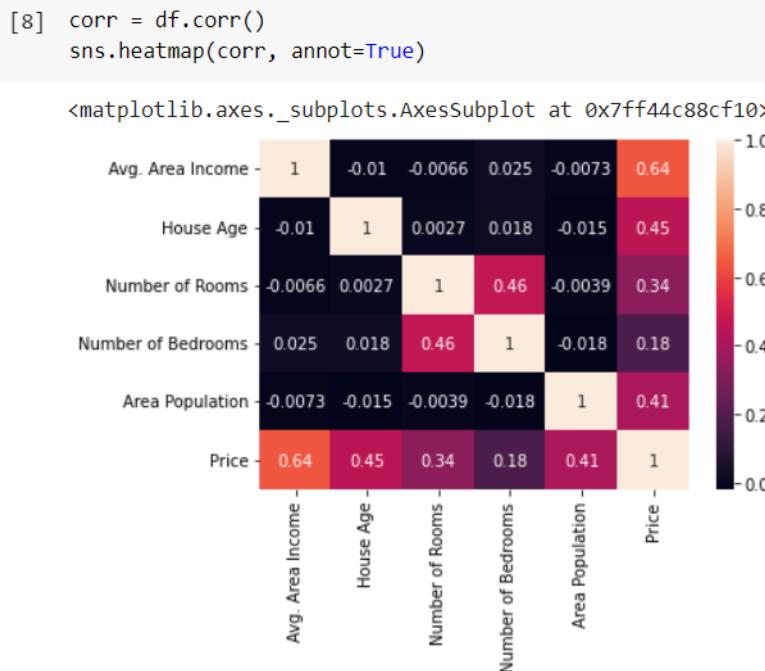


Once we've got a nice cleaned dataset, the next step is Exploratory Data Analysis (EDA). EDA is the process of figuring out what the data can tell us and we use EDA to find patterns, relationships, or anomalies to inform our subsequent analysis. While there are an almost overwhelming number of methods to use in EDA, one of the most effective starting tools is the pairs plot (also called a scatterplot matrix). A pairs plot allows us to see both distribution of single variables and relationships between two variables. Pair plots are a great method to identify trends for follow-up analysis.

One more step is to check out the correlation of the dataset by using the Python code:

```
1 corr = df.corr()  
2 sns.heatmap(corr, annot=True)
```

## Output



We now have come all the way of the processing and studying the data. This dataset as we have informed before, used to perform Regression model. The price that we wish predict is associated with other attributes. To perform Lasso Regression, we first need to separate the data to the Columns of X for input, and Columns of y for output. The price column is the sixth column of the dataset or we can say the last column, to take it away, we will do as follow:

```
1 data = df.values  
2 X, y = data[:, :-1], data[:, -1]
```

The training step is simple by applying the written function compare to the complicate of matrix multiply by Math base approach:

```
1 model = Lasso(alpha = 1.0)  
2 model.fit(X,y)
```

For the checking step, we can use the function `model.coef_` or `model.intercept_` to get the coefficient as well as the intercept of this linear model.

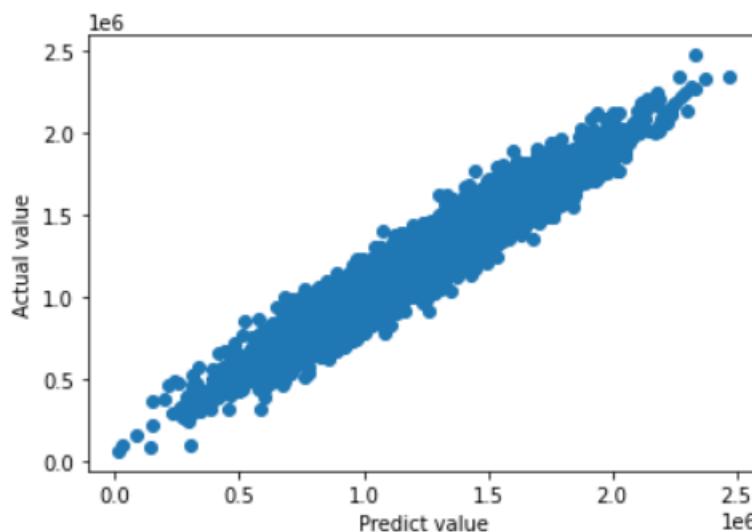
Another plot can be applied to check the result is the  $y = x$  scatter plot, actual data will be the y axis while the predict data is the x axis, if the scatter point lie perfectly on the  $y = x$  line. Our regression model predict the data have the low MSE.

```
1 actual = y
2 predict = model.predict(X)
3 plt.figure()
4 plt.scatter(actual, predict)
5 plt.xlabel('Predict value')
6 plt.ylabel('Actual value')
7 plt.show()
```

### Output



```
actual = y
predict = model.predict(X)
plt.figure()
plt.scatter(actual, predict)
plt.xlabel('Predict value')
plt.ylabel('Actual value')
plt.show()
```



The scatter points lie approximately to a line that is  $y = x$  show that our model predict quite good with the original data.



### 5.3 Implementation on web

We also program the implementation algorithm for the dataset on the website, here is its interface:

The screenshot shows a web-based application for estimating house prices. The interface includes the following fields and a button:

- Average of Income:** 79545.45857
- House age:** 5.682861322
- Number of rooms:** 7.009188143
- Number of bedrooms:** 4.09
- Area population:** 23086.8005
- Estimate Price:** (button)
- Result:** 1225032.5 USD

Problem demonstration on web

Source code of the web: [Github link](#)

## 6 Canada CO<sub>2</sub> Emissssions By Cars - Nam Trần

### 6.1 Data Overview

In recent years, the amount of gas in cars has been increased, which leads to serious air pollution in cities, especially megacities. That results in raising public concern about car use and has put pressure on car manufacturers to find a way to estimate the emissions their cars produce from some fundamental car features namely: model, class, engine. So in this section, we will try to build a machine learning model (Lasso Regression) to predict the car emissions. But before it, let's look at our dataset.

	MODEL	MAKE	VEHICLE CLASS	ENGINE_SIZE	CYLINDERS	TRANSMISSION	FUEL	FUEL_CONSUMPTION*	CO2_EMISSIONS
2	2001 ACURA	1.7EL	COMPACT	1.7	4 A4	X		9.3	191
3	2001 ACURA	1.7EL	COMPACT	1.7	4 M5	X		7.2	8.3
4	2001 ACURA	3.2CL	COMPACT	3.2	6 A55	Z		8.9	34
5	2001 ACURA	3.2TL	MID-SIZE	3.2	6 A55	Z		13.7	265
6	2001 ACURA	3.5RL	MID-SIZE	3.5	6 A4	Z		8.8	11.5
7	2001 ACURA	INTEGRA SUBCOMPACT	1.8	4 A4	X			13.8	24
8	2001 ACURA	INTEGRA SUBCOMPACT	1.8	4 M5	X			15	22
9	2001 ACURA	INTEGRA SUBCOMPACT	1.8	4 M5	Z			10.9	13.1
10	2001 ACURA	MDX	SUV	3.5	6 A5	Z		11.4	22
11	2001 ACURA	NSX	TWO-SEATER	3	6 A54	Z		15.5	301
12	2001 ACURA	NSX	TWO-SEATER	3.2	6 M6	Z		15.3	21
13	2001 AUDI	A4	COMPACT	1.8	4 A5	Z		15.6	306
14	2001 AUDI	A4	COMPACT	1.8	4 M5	Z		13.2	13.4
15	2001 AUDI	A4	COMPACT	2.8	6 A5	Z		8.1	25
16	2001 AUDI	A4 QUATT COMPACT	1.8	4 A5	Z			11.8	232
17	2001 AUDI	A4 QUATT COMPACT	1.8	4 M5	Z			14.6	285
18	2001 AUDI	A4 QUATT COMPACT	2.8	6 A5	Z			9.7	22
19	2001 AUDI	A4 QUATT COMPACT	2.8	6 M5	Z			13.3	265
20	2001 AUDI	A6	MID-SIZE	2.8	6 A5	Z		12.2	27
21	2001 AUDI	A6 AVANT STATION WAGON - MID-SIZE	2.8	6 A5	Z			10.1	244
22	2001 AUDI	A6 QUATT MID-SIZE	2.7	6 A5	Z			14.9	292
23	2001 AUDI	A6 QUATT MID-SIZE	2.7	6 M6	Z			14.3	288
24	2001 AUDI	A6 QUATT MID-SIZE	2.8	6 A5	Z			15	294
								10.5	306
								15.2	22
								10.7	301
								15.2	304
								10.7	21
								15.4	306

Canada Emission Car

Our dataset have total 11 comlumns (features) and 780 rows. However, we just need to focus on the following columns to make predict.

- *MAKE*: Car company which make the car.
- *MODEL*: Car model.
- *VEHICLE CLASS* : Class that the car belongs.
- *ENGINE\_SIZE*: The max size of the engine in cubic cm (cc).
- *CYLINDERS*: Number of cylinders.
- *TRANSMISSION*: Type of transmission.
- *FUEL*: Fuel Class.
- *TRANSMISSION*: Type of transmission.
- *FUEL\_CONSUMPTION*: Amount of Fuel consumption
- *CO<sub>2</sub>\_EMISSIONS*: Car emissions (g CO<sub>2</sub>/km). And it's also the label or taget of our later model.

That's quite enough to introduce about our car emissions dataset. Now, let's move into the main part that is "Implementaion" which describe detaily how and why the Lasso will be built.



## 6.2 Implementation

Before, we walk you through this most interested subsection, again we'll introduce the meta data or the catalog of it.

1. Preprocessing data: Importing libraries, importing datasets, clean missing data, encoding categorical, splitting dataset into training and test set, feature scaling
2. Model build: build the best Lasso model by sklearn library.
3. Evaluation the lasso model and compared with the other model such as: Linear Regression, Decision Tree, Ridge Regression.

### 6.2.1 Data Preprocessing

Data preprocessing is a process of preparing the raw data and making it suitable for a machine learning model. It is the first and crucial step while creating a machine learning model.

#### 6.2.1.a Import Library

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import matplotlib.lines as mlines
4 import matplotlib.transforms as mtransforms
5 from sklearn.linear_model import LinearRegression
6 from sklearn.model_selection import ShuffleSplit
7 from sklearn.model_selection import cross_val_score
8 from sklearn.model_selection import GridSearchCV
9 from sklearn.linear_model import Lasso
10 from sklearn.tree import DecisionTreeRegressor
11 import pandas as pd
12 import seaborn as sns
```

All the above libraries are the fundamental libraries, we'll use to finish the Implementation segment. For example: numpy library to efficiently working with matrix, sklearn Lasso Regression, Linear Regression, DecisionTree Regression model.

#### 6.2.1.b Import Data

```
1 df = pd.read_csv("Sample.csv")
2 df.head()
```

We assign the dataset into a dataframe called [df] then we check to see that whether the data is successful import or not.

Output:



	MODEL	MAKE	MODEL.1	...	Unnamed: 10	Unnamed: 11	CO2_EMISSIONS
0	2001	ACURA	1.7EL	...	8.3	34	191
1	2001	ACURA	1.7EL	...	8.3	34	191
2	2001	ACURA	3.2CL	...	11.5	25	265
3	2001	ACURA	3.2TL	...	11.6	24	267
4	2001	ACURA	3.5RL	...	13.1	22	301

[5 rows x 13 columns]

Now, we will drop some columns and just keep the important columns as we recommended at the beginning.

```
1 df2 = df.drop(['MODEL', 'Unnamed: 9', 'Unnamed: 10', 'Unnamed: 11'], axis = 1)
2 df2.head()
```

	MAKE	MODEL.1	VEHICLE CLASS	ENGINE_SIZE	CYLINDERS	TRANSMISSION	FUEL	FUEL_CONSUMPTION*	CO2_EMISSIONS
0	ACURA	1.7EL	COMPACT	1.7	4	A4	X	9.3	191
1	ACURA	1.7EL	COMPACT	1.7	4	M5	X	8.9	191
2	ACURA	3.2CL	COMPACT	3.2	6	AS5	Z	13.7	265
3	ACURA	3.2TL	MID-SIZE	3.2	6	AS5	Z	13.8	267
4	ACURA	3.5RL	MID-SIZE	3.5	6	A4	Z	15.0	301

Data After Drop

### 6.2.1.c Clean Missing Data

Here, we'll find the rows which contain the values "NA" then remove it. First, let's check how many missing data each columns has.

```
1 df2.isnull().sum()
```

```
MAKE          0
MODEL.1       0
VEHICLE CLASS 0
ENGINE_SIZE    0
CYLINDERS      0
TRANSMISSION    0
FUEL           0
FUEL_CONSUMPTION* 0
CO2_EMISSIONS   0
dtype: int64
```

Check Null Columns

Luckily, we do not have any null values in each column. So nothing to do more in this step. Let's move on the next step.



#### 6.2.1.d Encoding Categorical

Although our dataset do not have any null values, the dataset contains until five columns which contains categorical data. And obviously, the Lasso can not work with the text data, it must to convert to corresponding label. So, in this steps, we'll label all the categorical data in each column to make it satisfied the Lasso model input.

Columns **MAKE**:

First, we'll check how many unique values in the column and label it by an integer number. Let's start

```
[1] df2['MAKE'].unique()
```

```
array(['ACURA', 'AUDI', 'BMW', 'BUICK', 'CADILLAC', 'CHEVROLET',
       'CHRYSLER', 'DAEWOO', 'DODGE', 'FORD', 'GMC', 'HONDA', 'HYUNDAI',
       'INFINITI', 'ISUZU', 'JAGUAR', 'JEEP', 'KIA', 'LEXUS', 'LINCOLN',
       'MAZDA', 'MERCEDES-BENZ', 'NISSAN', 'OLDSMOBILE', 'PLYMOUTH',
       'PONTIAC', 'PORSCHE', 'SAAB', 'SATURN', 'SUBARU', 'SUZUKI',
       'TOYOTA', 'VOLKSWAGEN', 'VOLVO'], dtype=object)
```

Before Label Column MAKE

So, we have exactly 34 distinguished values. And it's small enough to label it by hand instead using **one hot encoding**. In the other words, the first distinguished value will be labeled as 0, then next distinguished value is 1,...until 33.

At the end, we have a new MAKE column with the numerical content, which follow the rules as above.

```
array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18,
       19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33],
      dtype=object)
```

After Label Column MAKE

Similar encoding rule as column **MAKE**, with three left columns namely: TRANSMISSION, VEHICLE CLASS, FUEL, we do same things.

Columns **VEHICLE CLASS**:

Before encoding:

```
[13] df2['VEHICLE CLASS'].unique()

array(['COMPACT', 'MID-SIZE', 'SUBCOMPACT', 'SUV', 'TWO-SEATER',
       'STATION WAGON - MID-SIZE', 'FULL-SIZE', 'MINICOMPACT',
       'STATION WAGON - SMALL', 'VAN - CARGO', 'VAN - PASSENGER',
       'PICKUP TRUCK - STANDARD', 'PICKUP TRUCK - SMALL', 'MINIVAN'],
      dtype=object)
```

Before Encode Column VEHICLE CLASS



After encoding:

```
df2['VEHICLE CLASS'].unique()  
array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13], dtype=object)
```

After Encode Column VEHICLE CLASS

Columns **TRANSMISSION**:

Before encoding

```
df2['TRANSMISSION'].unique()  
array(['A4', 'M5', 'AS5', 'A5', 'AS4', 'M6', 'A3', 'AV'], dtype=object)
```

Before Encode Column TRANSMISSION

After encoding:

```
df2['TRANSMISSION'].unique()  
array([0, 1, 2, 3, 4, 5, 6, 7], dtype=object)
```

After Encode Column TRANSMISSION

Columns **FUEL**:

Before encoding

```
[ ] df2['FUEL'].unique()  
array(['X', 'Z', 'E', 'N', 'D'], dtype=object)
```

Before Encode Column FUEL

After encoding:



```
df2['FUEL'].unique()

array([0, 1, 2, 3, 4], dtype=object)
```

After Encode Column FUEL

The first four columns can be applied by the handing label, but the last column that is **MODEL.1** can not be address as same way. Because, the number of unique values of it is very huge. There

```
df2['MODEL.1'].unique()

array(['1.7EL', '3.2CL', '3.2TL', '3.5RL', 'INTEGRA', 'INTEGRA GS-R',
       'MDX', 'NSX', 'AA', 'A QUATTRO', 'A6', 'A6 AVANT QUATTRO WAGON',
       'A6 QUATTRO', 'AB L', 'AB QUATTRO', 'ALLROAD', 'S4 QUATTRO',
       'S8 QUATTRO #', 'TT COUPE QUATTRO', 'TT COUPE QUATTRO #',
       'TT ROADSTER', 'TT ROADSTER QUATTRO #', 'TTS', 'TTSI',
       '325Ci CONVERTIBLE', '325i', '325xi', '330Ci', '330i', '330i CONVERTIBLE',
       '330i', '330xi', '525i', '525i SPORT WAGON', '530i', '540i',
       '540i SPORT WAGON', '740i', '740iL', '750iL', 'COUPE',
       'M ROADSTER', 'M3', 'M3 CONVERTIBLE', 'MS', 'XS', 'Z3', 'Z8',
       'CENTURY', 'LESABRE', 'PARK AVENUE', 'PARK AVENUE #', 'REGAL',
       'REGENCY', 'CATERAT', 'DEVILLE', 'ELDORADO', 'SEVILLE',
       'ASTRO CARGO', 'ASTRO VAN AND', 'ASTRO PASSENGER',
       'ASTRO PASSENGER AND', 'BLAZER', 'BLAZER 4X4', 'C1500 SILVERADO',
       'C1500 SUBURBAN', 'C1500 TAHOE', 'Camaro', 'CAVALIER',
       'CAVALIER BI-FUEL', 'CORVETTE', 'G15/G20 CHEVY VAN',
       'G1500/G2500 CHEVY EXPRESS', 'IMPALA', 'K1500 SILVERADO 4X4',
       'K1500 SUBURBAN 4X4', 'K1500 TAHOE 4X4', 'MONTIBU', 'MONTE CARLO',
       'S10 FFV', 'S10 4X4', 'S10 FFV 4X4', 'SEBRING CONVERTIBLE 4X4',
       'TRACKER VAN 4X4', 'VENTURE', '300M', 'CONCORD', 'INTREPID',
       'LHS', 'NEON', 'NEON R/T #', 'PT CRUISER', 'SEBRING',
       'SEBRING CONVERTIBLE', 'SEBRING COUPE', 'TOWN & COUNTRY',
       'TOWN & COUNTRY AWD', 'LANOS', 'LEGANZA', 'NUBIRA', 'NUBIRA WAGON',
       'CARAVAN', 'DAKOTA', 'DAKOTA 4X4', 'DURANGO 4X4',
       'DURANGO 4X4', 'F150', 'F150 4X4', 'GRAND CARAVAN', 'CARAVAN AWD',
       'RAM 1500', 'RAM 1500 4X4', 'RAM 1500 VAN', 'RAM 1500 WAGON',
       'RAM 2500 VAN', 'RAM 2500 VAN CNG', 'RAM 2500 WAGON',
       'RAM 2500 WAGON CNG', 'VIPER GTS', 'VIPER RT/10', 'COUGAR',
       'CROWN VICTORIA', 'CROWN VICTORIA NGV', 'E150 CLUB WAGON',
       'E150 VAN', 'E250 VAN', 'ESCAPE', 'ESCAPE 4X4',
```

Before Group FUEL Columns

are 351 unique values in this column, so it will be cause a lot of variance if we label it as the previous way. Instead, we use **one-hot encoding** to encode this columns. But first, we will group all the uniques value which has a frequency less than 3 in the dataset and name it as 'other'.

```
1 Model_stats = df2.groupby('MODEL.1')['MODEL.1'].agg('count').sort_values(ascending
   =False)
2 Model_stats_less_than_3 = Model_stats[Model_stats <= 3]
3 df2['MODEL.1'] = df2['MODEL.1'].apply(lambda x: 'other' if x in
   Model_stats_less_than_3 else x)
4 df2['MODEL.1'].unique()
```

Output:

```
array(['other', 'A QUATTRO', 'A6 QUATTRO', 'Z3', 'C1500 SILVERADO',
       'CAMARO', 'CAVALIER', 'K1500 SILVERADO 4X4', 'S10 FFV', 'INTREPID',
       'SEBRING COUPE', 'LANOS', 'DAKOTA', 'DAKOTA 4X4', 'RAM 1500',
       'F150', 'F150 4X4', 'MUSTANG', 'RANGER', 'RANGER 4X4',
       'C1500 SIERRA', 'K1500 SIERRA 4X4', 'SONOMA FFV', 'ACCENT',
       'TJ 4X4', 'MAGENTIS', 'LS', '626', 'PROTEGE', 'SENTRA', 'FIREBIRD',
       'SUNFIRE', 'IMPREZA AWD', 'CELICA', 'TACOMA', 'TACOMA 4X4',
       'TUNDRA', 'GOLF', 'JETTA', 'NEW BEETLE', 'PASSAT', 'PASSAT WAGON'],
      dtype=object)
```

Now, the number of unique values has greatly decreased but still large (50 distinguished values). So, we'll use one hot encoding at this moment.

```
1 dummies = pd.get_dummies(df2['MODEL.1'])
```



```
2 df2 = pd.concat([df2.drop('MODEL_1', axis = 1), dummies.drop('other',axis='columns')],axis='columns')
3 df2.head(3)
```

Output:

MAKE	VEHICLE_CLASS	ENGINE_SIZE	CYLINDERS	TRANSMISSION	FUEL_FUEL_CONSUMPTION*	CO2_EMISSIONS	626	A4	A6	ACCENT	C1500	C1500	CAVALIER	CAVALIER	CELICA	DAKOTA	DAKOTA	F150	F150	FIREBIRD	FIREBIRD
0	0	1.7	4	0 0	9.3	191	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	1.7	4	1 0	8.9	191	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	3.2	6	2 1	13.7	265	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

The number of columns in our dataframe has increased result in expanding the dataframe size. However, due to the small samples in the dataset (more than 780 samples) so it is acceptable.

### 6.2.1.e Data Split

As many machine learning model, we'll assign all the input features into **X** and the label in **y**. Then we'll split the training set and testing set by the ratio 4/5 (20% for testing, 80% for training).

```
1 X_train, X_test, y_train, y_test = train_test_split(X,y,test_size=0.2,random_state
=10)
```

### 6.2.2 Build Lasso Model

Now, we'll build a Lasso Regression basing on the library **sklearn.linear\_model import Lasso** by some lines code.

```
1 Lasso_model = Lasso(alpha = 1)
2 Lasso_model.fit(X_train, y_train)
```

After that, since we use **one-hot encoding** so we can not directly use the predict method from sklearn, we must modify the data input into standard form by declare a new function called **predict\_emission**.

```
1 def predict_emission(Company ,VEHICLE_CLASS ,ENGINE_SIZE ,CYLINDERS , TRANSMISSION ,
FUEL , FUEL_CONSUMPTION , Model):
2     loc_index = np.where(X.columns==Model)[0][0]
3     x = np.zeros(len(X.columns))
4     x[0] = Company
5     x[1] = VEHICLE_CLASS
6     x[2] = ENGINE_SIZE
7     x[3] = CYLINDERS
8     x[4] = TRANSMISSION
9     x[5] = FUEL
10    x[6] = FUEL_CONSUMPTION
11    if loc_index >= 0:
12        x[loc_index] = 1
13    return Lasso_model.predict([x])[0]
14 predict_emission(0, 2, 3, 3, 2,1 ,22, 'RANGER')
```

Output:



426.76853869080946

### 6.2.3 Evaluation

After building the Lasso model, we should evaluate to see whether this model did a good job with our dataset or not. Then, we'll compare the lasso regression with other regression such as Linear and Decision regression when apply in this dataset.

#### 6.2.3.a Lasso Evaluation

As we splitted our dataset into two main part (testing and training), now we are going to use the testing set to evaluate the model. Furthermore, we use a method called "cross validation score" to evaluate model.

```
1 from sklearn.model_selection import ShuffleSplit
2 from sklearn.model_selection import cross_val_score
3 cv = ShuffleSplit(n_splits = 5, test_size = 0.2, random_state = 0)
4 cross_val_score(Lasso(alpha = 1), X, y, cv = cv)
```

Output:

```
array([0.90757461, 0.89229103, 0.93831544, 0.90934974, 0.86988646])
```

As we can see that, almost the case, our model get nearly 90 percent score, which are very high. This also implies that our model actually did a good job on such car emission dataset.

### 6.2.4 Models Comparasion

To make everything simple, we still use cross validation as a way to evaluate three following models: Lasso, Linear, and Decision Tree Regression. Then compare these three model with each other.

```
1 def find_best_model_using_gridsearchcv(X,y):
2     algos = {
3         'linear_regression': {
4             'model': LinearRegression(),
5             'params':{
6                 'normalize': [True, False]
7             }
8         },
9         'lasso': {
10             'model': Lasso(),
11             'params': {
12                 'alpha': [1,2],
13                 'selection': ['random', 'cyclic']
14             }
15         },
16         'decision_tree': {
17             'model': DecisionTreeRegressor(),
18             'params': {
```



```
19         'criterion': ['mse', 'friedman_mse'],
20         'splitter': ['best', 'random']
21     }
22 }
23 scores = []
24 cv = ShuffleSplit(n_splits=5, test_size = 0.2, random_state = 0)
25 for algo_name, config in algos.items():
26     gs = GridSearchCV(config['model'], config['params'], cv=cv,
27     return_train_score = False)
28     gs.fit(X,y)
29     scores.append({
30         'model': algo_name,
31         'best_score': gs.best_score_,
32         'best_params': gs.best_params_
33     })
34 return pd.DataFrame(scores, columns=['model','best_score','best_params'])
35
36 find_best_model_using_gridsearchcv(X,y)
```

Output:

	model	best_score	best_params
0	linear_regression	0.920526	{'normalize': True}
1	lasso	0.903483	{'alpha': 1, 'selection': 'cyclic'}
2	decision_tree	0.948663	{'criterion': 'friedman_mse', 'splitter': 'best'}

The function above do two main tasks.

- Firstly, it chooses the best tuning parameter to each algorithm.
- Secondly, it print all best score of each algorithm.

As we can see, the Decision Tree Regression has the highest score or it is the most suit algorithm on this dataset, following by the Linear Regression then the last one is Lasso. But, in terms of regularization, Lasso can ensure that, the features do not overweights each other and give the real life input estimation as near as possible when training.



## 7 COVID-19 vaccination vs mortality - Quốc Việt

### 7.1 Data Overview

The COVID-19 pandemic has brought the whole planet to its knees. More over 5 million people have died since the writing of this report, and the only acceptable way out of the disaster is to vaccinate all parts of society.

In spite of the fact that the benefits of vaccination have been proved to the world many times, anti-vaccine groups are springing up all over the world. This is the data set which aim to investigate the impact of coronavirus vaccinations on coronavirus mortality.

We are going to take a overview about this data set. Since it is time when each country try to vaccinate people as much as possible, this is so perfect for a regression model in general or just only the idea for LASSO. The data contains the following columns:

- **Country:** country name
- **ISO\_code:** iso code for each country
- **Date:** The date the data record
- **Total\_vaccinations:** number of all doses of COVID vaccine usage in that country
- **People\_vaccinated:** number of people who got at least one shot of COVID vaccine
- **People\_fully\_vaccinated:** number of people who got full vaccine shots
- **New\_deaths:** number of daily new deaths
- **Population:** 2021 country population
- **Ratio:** percentage of vaccinations in that country at that date =  $\text{people\_vaccinated} / \text{population}$



## 7.2 Running with Python code

The first step is to import some library that involve in the demonstration.

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import pandas as pd
4 import seaborn as sns
5 from sklearn.linear_model import Lasso
```

The library **numpy** is a Python library used for working with arrays. It also has functions for working in domain of linear algebra, fourier transform, and matrices.

The part of **matplotlib**, matplotlib.pyplot is a collection of functions that make matplotlib work like MATLAB. Each pyplot function makes some change to a figure: e.g., creates a figure, creates a plotting area in a figure, plots some lines in a plotting area, decorates the plot with labels, etc. The library **pandas** is a fast, powerful, flexible and easy to use open source data analysis and manipulation tool, built on top of the Python programming language.

The library **seaborn** is a Python data visualization library based on matplotlib. It provides a high-level interface for drawing attractive and informative statistical graphics.

Finally, the tool for our model is imported from **sklearn** - a Python module for machine learning built on top of SciPy and is distributed under the 3-Clause BSD license, with the **linear\_model** of **Lasso**.

To import the data, we will follow the id of the link on Google Drive. This makes the dataset accessible by many people in order to view it, just need to follow the ID.

```
1 %%shell
2 covid="1X70Y36vyn3k60mE7mG4VRs8oza3yGbOT"
3
4 cd /content/
5 gdown -q --id $covid
```

Next, we check if the data is imported successfully, if Yes then print 5 rows in the head and 5 rows in the tail, check it with our own dataset to ensure that it is correct.

```
1 df = pd.read_csv("covid_vaccination_vs_death_ratio.csv")
2 print(df.head(), df.tail())
```

### Output

```
df = pd.read_csv("covid_vaccination_vs_death_ratio.csv")
print(df.head(), df.tail())

   Unnamed: 0      country iso_code ... New_deaths population      ratio
0        0  Afghanistan     AFG ...       12  40094444.0  1.119552
1        1  Afghanistan     AFG ...       10  40094444.0  1.173083
2        2  Afghanistan     AFG ...       10  40094444.0  1.188112
3        3  Afghanistan     AFG ...       19  40094444.0  1.195607
4        4  Afghanistan     AFG ...       14  40094444.0  1.196111

[5 rows x 10 columns]   Unnamed: 0      country iso_code ... New_deaths population      ratio
20327      20327  Zimbabwe     ZWE ...       2  15158323.0  20.566893
20328      20328  Zimbabwe     ZWE ...       0  15158323.0  20.639473
20329      20329  Zimbabwe     ZWE ...       0  15158323.0  20.717239
20330      20330  Zimbabwe     ZWE ...       3  15158323.0  20.797927
20331      20331  Zimbabwe     ZWE ...       3  15158323.0  20.858567

[5 rows x 10 columns]
```



Next, we will get deeper view into the data by using the `df.info()`, this function will return all the columns of the dataset together with the Data type (Dtype).

### Output

```
df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 20332 entries, 0 to 20331
Data columns (total 10 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   Unnamed: 0        20332 non-null   int64  
 1   country          20332 non-null   object  
 2   iso_code         20332 non-null   object  
 3   date             20332 non-null   object  
 4   total_vaccinations 20332 non-null   float64 
 5   people_vaccinated 20332 non-null   float64 
 6   people_fully_vaccinated 20332 non-null   float64 
 7   New_deaths       20332 non-null   int64  
 8   population        20332 non-null   float64 
 9   ratio             20332 non-null   float64 
dtypes: float64(5), int64(2), object(3)
memory usage: 1.6+ MB
```

With the data set we have here, it is so perfect that all the rows in the Dataset is non-null show that a good data set.. This results also in the benefits that we don't have to perform some changes to our data to best fit the model. The Data cleaning process is also covered by the non-null, we can also check the number of null value by using another function of **pandas: df.isnull()**.

The Selection of attributes step is performed after we get the description about the dataset with the supported function `df.describe()`:

### Output

```
[28] df.describe()
```

	Unnamed: 0	total_vaccinations	people_vaccinated	people_fully_vaccinated	New_deaths	population	ratio
count	20332.000000	2.033200e+04	2.033200e+04	2.033200e+04	20332.000000	2.033200e+04	20332.000000
mean	10165.500000	1.798981e+07	1.155247e+07	6.571098e+06	108.781871	4.863896e+07	29.920624
std	5869.487172	6.328884e+07	4.242151e+07	2.284904e+07	353.037676	1.559664e+08	24.936193
min	0.000000	1.630000e+02	1.620000e+02	1.000000e+00	-41.000000	1.619000e+03	0.003023
25%	5082.750000	4.106325e+05	2.971550e+05	9.506200e+04	1.000000	2.871380e+06	7.197272
50%	10165.500000	2.363600e+06	1.548864e+06	7.438710e+05	10.000000	1.015944e+07	24.266101
75%	15248.250000	1.059296e+07	6.060886e+06	4.137016e+06	58.000000	3.816076e+07	50.225309
max	20331.000000	2.174043e+09	1.100842e+09	1.022207e+09	8786.000000	1.445585e+09	118.567868



Next is the preprocessing step, we will eliminate the null values and check the columns of this data set:

```
▶ Total_null_values = df.isnull().sum().sum()
print(Total_null_values)

0

[8] df.columns

Index(['Unnamed: 0', 'country', 'iso_code', 'date', 'total_vaccinations',
       'people_vaccinated', 'people_fully_vaccinated', 'New_deaths',
       'population', 'ratio'],
      dtype='object')
```

There is a column namely "**Unnamed**" as it is the column for indices, we will need to eliminate it with the code:

```
▶ for col in df.columns:
    if col.startswith('Unnamed'):
        del df[col]

df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 20332 entries, 0 to 20331
Data columns (total 9 columns):
 #   Column            Non-Null Count  Dtype  
--- 
 0   country           20332 non-null   object  
 1   iso_code          20332 non-null   object  
 2   date              20332 non-null   object  
 3   total_vaccinations 20332 non-null   float64 
 4   people_vaccinated 20332 non-null   float64 
 5   people_fully_vaccinated 20332 non-null   float64 
 6   New_deaths        20332 non-null   int64   
 7   population        20332 non-null   float64 
 8   ratio              20332 non-null   float64 
dtypes: float64(5), int64(1), object(3)
memory usage: 1.4+ MB
```

After checking for the columns, because this data set is about Covid-19, we need to check the country that give the data and plot some figures about it.

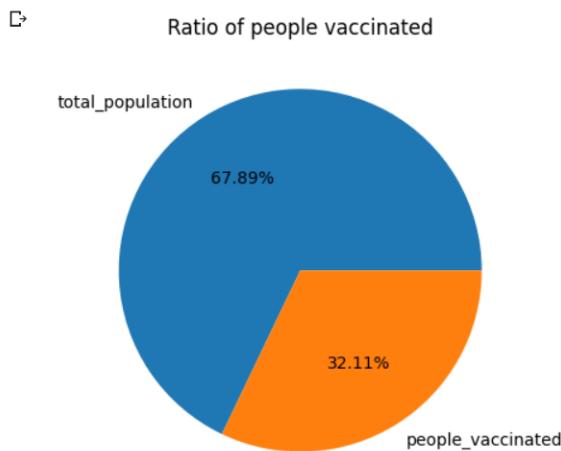
### Output

```
[10] # The number of countries
df['country'].unique().shape[0]
```

This is the code for plotting the percentage of one-shot vaccinated people over the total population.

```
1 # The total population of all countries recorded
2 total_pop = df['population'].unique().sum()
3 print(total_pop)
4 # Checking the number of vaccinated people
5 group_country_1 = df.groupby('country').agg({'people_vaccinated':max})
6 vaccine_1shot = group_country_1['people_vaccinated'].sum()
7 print(vaccine_1shot)
8 # Checking the number of fully vaccinated people
9 group_country_2 = df.groupby('country').agg({'people_fully_vaccinated':max})
10 vaccine_2shot = group_country_2['people_fully_vaccinated'].sum()
11 print(vaccine_2shot)
12
13 plt.figure(figsize=(10,5),dpi=100)
14 plt.pie(
15     [total_pop, vaccine_1shot],
16     autopct='%.2F%%',
17     labels=['total_population','people_vaccinated'])
18 plt.title('Ratio of people fully vaccinated')
19 plt.show()
```

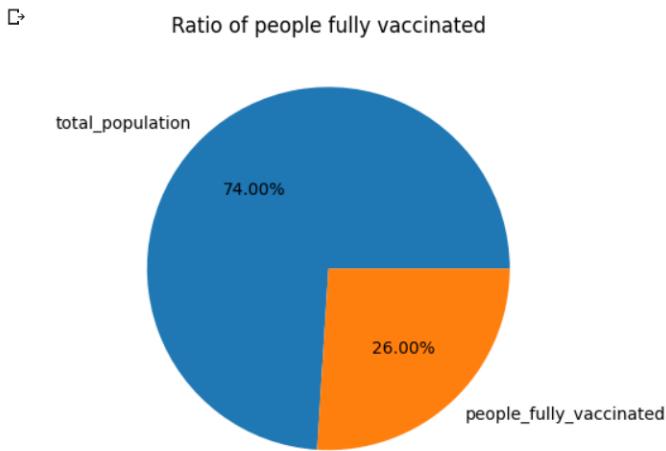
## Output



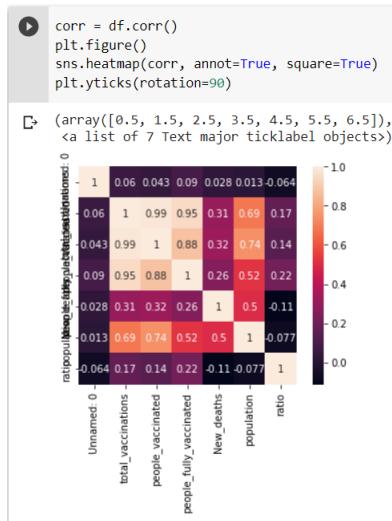
This is the code for plotting the percentage of fully vaccinated people over the total population.

```
1 plt.figure(figsize=(10,5),dpi=100)
2 plt.pie(
3     [total_pop, vaccine_2shot],
4     autopct='%.2F%%',
5     labels=['total_population','people_fully_vaccinated'])
6 plt.title('Ratio of people fully vaccinated')
7 plt.show()
```

## Output



And the heatmap from Seaborn for the correlation matrix:

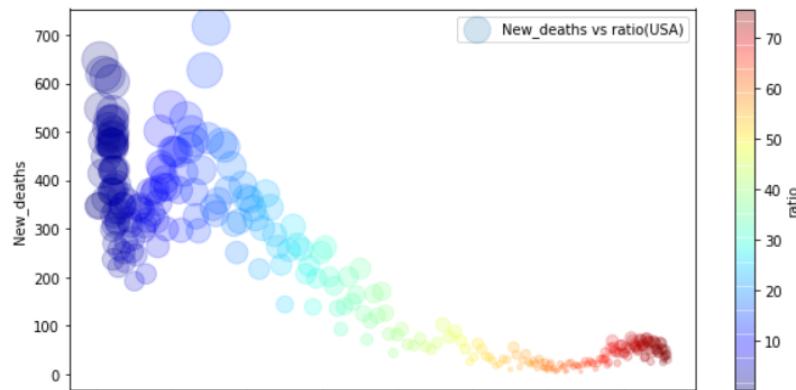


### 7.3 Fitting model - Covid-19 Italy

Italy has administered at least 90,830,635 doses of COVID vaccines the time that this report been written. Assuming every person needs 2 doses, that's enough to have vaccinated about 75.3 percent of the country's population.

Let's checkout the plotting and model fitting. Starting with plotting the data of New deaths like the previous section about the USA.

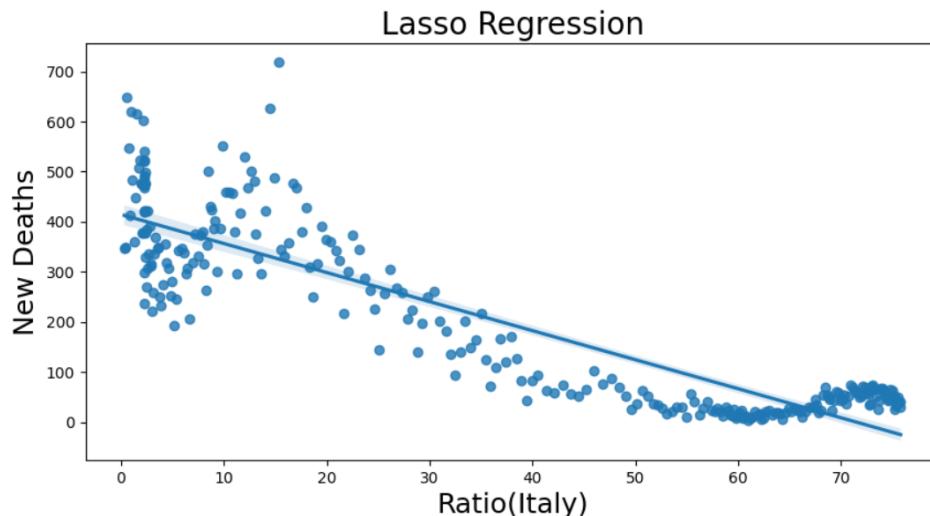
```
1 df_italy = df[df['country']=='Italy']
2 df_italy.plot(
3     kind='scatter',
4     x='ratio',
5     y='New_deaths',
6     label='New_deaths vs ratio(USA)',
7     alpha=0.2,
8     figsize=(10,5),
9     s=df_italy['New_deaths'],
10    c='ratio',
11    cmap=plt.get_cmap('jet'),
12    colorbar=True)
```



Fitting the model by LASSO with take the code below:

```
1 # define x,y
2 x_italy = df_italy[['ratio']]
3 y_italy = df_italy[['New_deaths']]
4
5 # fitting in
6 lasso_italy = Lasso(alpha=1.0)
7 lasso_italy.fit(x_italy,y_italy)
8
9 # coefficients
10 print ('Coefficients: ', lasso_italy.coef_)
11 print ('Intercept: ', lasso_italy.intercept_)
12
13 plt.figure(figsize=(10,5),dpi= 100)
14 sns.regplot(x= x_italy, y= y_italy)
15 plt.xlabel('Ratio(Italy)',fontsize=18)
16 plt.ylabel('New Deaths',fontsize=18)
17 plt.title('Lasso Regression',fontsize=20)
```

## Output



## 7.4 Evaluation

After building the Lasso model, we should evaluate to see whether this model did a good a job with our dataset or not. Then, we'll compare the lasso regression with other regression such as Linear and Decision regression when apply in this dataset.

As we splitted our dataset into two main part (testing and training), now we are going to use the testing set to evaluate the model. Furthermore, we use a method called "cross validation score" to evaluate model and then create a new function to compare to Linear Regression:

```
1 from sklearn.linear_model import LinearRegression
2 from sklearn.model_selection import ShuffleSplit
3 from sklearn.model_selection import cross_val_score
4 from sklearn.model_selection import GridSearchCV
5
6
```



```
7 def find_best_model_using_gridsearchcv(X,y):
8     algos = {
9         'linear_regression': {
10             'model': LinearRegression(),
11             'params':{
12                 'normalize': [True, False]
13             }
14         },
15         'lasso': {
16             'model': Lasso(),
17             'params': {
18                 'alpha': [1,2],
19                 'selection': ['random', 'cyclic']
20             }
21         }
22     }
23     scores = []
24     cv = ShuffleSplit(n_splits=5, test_size = 0.2, random_state = 0)
25     for algo_name, config in algos.items():
26         gs = GridSearchCV(config['model'], config['params'], cv=cv,
27         return_train_score = False)
28         gs.fit(X,y)
29         scores.append({
30             'model': algo_name,
31             'best_score': gs.best_score_,
32             'best_params': gs.best_params_
33         })
34     return pd.DataFrame(scores, columns=['model','best_score','best_params'])
```

Then apply the functions:

```
1 cv = ShuffleSplit(n_splits = 5, test_size = 0.2, random_state = 0)
2 cross_val_score(Lasso(alpha = 1), x_italy, y_italy, cv = cv)
```

**Output:**

```
cv = ShuffleSplit(n_splits = 5, test_size = 0.2, random_state = 0)
cross_val_score(Lasso(alpha = 1), x_italy, y_italy, cv = cv)

array([0.82243262, 0.80092539, 0.76284957, 0.78810394, 0.74180319])
```

Also when compare to normal Linear Regression:

```
1 find_best_model_using_gridsearchcv(x_italy, y_italy)
```

**Output:**

	model	best_score	best_params
0	linear_regression	0.783211	{'normalize': True}
1	lasso	0.783235	{'alpha': 2, 'selection': 'random'}

Lasso Regression has the benefit to reduce overfitting, when encounter to new data Lasso tend to perform better in some cases. In this, the prediction we get here show that Lasso slightly better than Linear Regression.



## 7.5 Conclusion

The Covid-19 pandemic make the world suffer a lot. From two countries of the data set we can tell that the higher the vaccinate ratio is, the decrement of new deaths will be. The data set also show a little strange data happen in Italy that the death ratio increase but with the period of time also in the data set we can conclude that:

- The number of new deaths raise due the time that people have vaccinated. Some people have get fully vaccinated in first of 2021 or late 2020, with the time when by, the dosages of vaccine in their body decrease - confirm in some researchs. [Check here](#)
- The new variant of Covid-19 namely Delta emerge. [Check here](#)



## 8 Factors influencing Life Expectancy (WHO) - Tân Phước

### 8.1 Inspiration and Data Overview

Improving life expectancy has always been a great concern of humanity for centuries. Hence, there have been lots of studies undertaken in the past on factors affecting life expectancy considering demographic variables, income composition and mortality rates. However, it was found that affect of immunization and human development index was not taken into account in the past. Also, some of the past research was done considering multiple linear regression based on data set of one year for all the countries.

Therefore, in this dataset, the Global Health Observatory (GHO) data repository under World Health Organization (WHO) keeps track of the health status as well as many other related factors for all countries. The data-set related to life expectancy, health factors for 193 countries has been collected from the same WHO data repository website and its corresponding economic data was collected from United Nation website. Among all categories of health-related factors only those critical factors were chosen which are more representative. It has been observed that in the past 15 years, there has been a huge development in health sector resulting in improvement of human mortality rates especially in the developing nations in comparison to the past 30 years. With that, WHO has considered data from year 2000-2015 for 193 countries for this dataset.

See the dataset: [Life expectancy dataset](#)

This dataset consists of 22 columns and 2938 rows. However, we just need to focus on the following columns to make predict:

- Status: Developed or Developing country
- Life expectancy: Life Expectancy in age
- Alcohol: Alcohol, recorded per capita (15+) consumption (in litres of pure alcohol)
- BMI: Average Body Mass Index of entire population
- Polio: Polio (Pol3) immunization coverage among 1-year-olds (%)
- HIV/AIDS: Deaths per 1 000 live births HIV/AIDS (0-4 years)
- GDP: Gross Domestic Product per capita (in USD)
- Income composition of resources: Human Development Index in terms of income composition of resources (index ranging from 0 to 1)
- Schooling: Number of years of Schooling(years)

Using these features, we will build and apply our Lasso Regression model to predict Life expectancy. Now, let's move to the main part.

### 8.2 Data analysis

Before jumping directly into building a model, we will first have to analyze and preprocess the data in this section. This step is vital to improve the performance and get a good result at the end.



### 8.2.1 Import library and data

To start with, we need to import some libraries that we are going to work with during the project. The followings are all of the necessary libraries:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import matplotlib.lines as mlines
4 import matplotlib.transforms as mtransforms
5 import pandas as pd
6 import seaborn as sns
7 from sklearn.linear_model import LinearRegression, Lasso
8 from sklearn.ensemble import RandomForestRegressor
9 from sklearn.model_selection import ShuffleSplit
10 from sklearn.model_selection import cross_val_score
11 from sklearn.model_selection import train_test_split
12 from sklearn import metrics
```

Brief information about these libraries:

- Numpy: is a Python library used for working with arrays. It also has functions for working in domain of linear algebra, fourier transform, and matrices. NumPy stands for Numerical Python.
- Matplotlib: is a comprehensive library for creating static, animated, and interactive visualizations in Python.
- Pandas: is a Python package that provides fast, flexible, and expressive data structures designed to make working with "relational" or "labeled" data easier and more intuitive.
- Seaborn: is a Python data visualization library based on matplotlib. It provides a high-level interface for drawing attractive and informative statistical graphics.
- 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, random forests, gradient boosting, k-means and DBSCAN, and is designed to interoperate with the Python numerical and scientific libraries NumPy and SciPy. In this project, we will use it to build our Lasso Regression model.

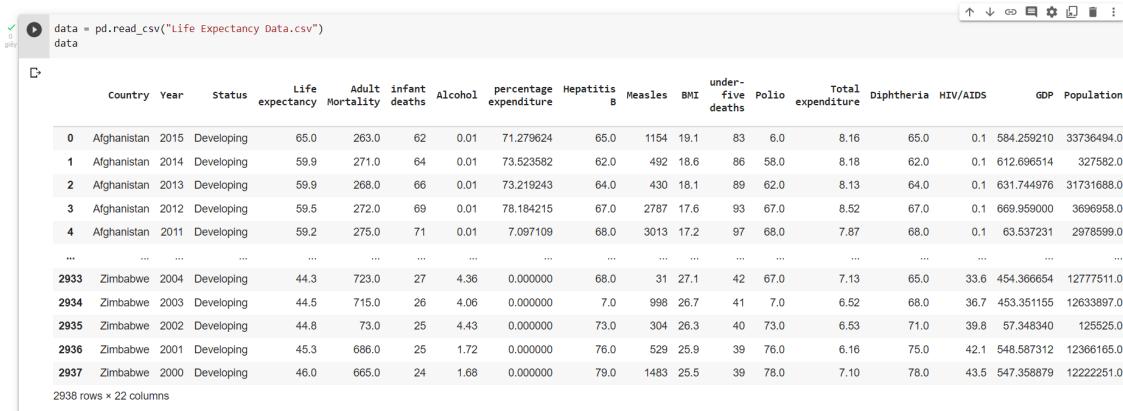
Now, we are going to import the dataset. To do this, we use the id of the link on Google Drive. This makes the dataset accessible to all people. See the following:

```
1 %%shell
2 cd /content
3 gdown -q --id "1iZ6BAL1VBLEfQPyuGL_NUf91rNyAoriI"
```

After importing the dataset successfully, we are now able to read everything from the dataset by running the code below:

```
1 data = pd.read_csv("Life Expectancy Data.csv")
2 data
```

**Output:**



A screenshot of a Jupyter Notebook cell. The code `data = pd.read\_csv("Life Expectancy Data.csv")` has been run, and the resulting DataFrame is displayed. The DataFrame has 2938 rows and 22 columns. The columns are: Country, Year, Status, Life expectancy, Adult Mortality, infant deaths, Alcohol, percentage expenditure, Hepatitis B, Measles, BMI, under-five deaths, Polio, Total expenditure, Diphtheria, HIV/AIDS, GDP, and Population.

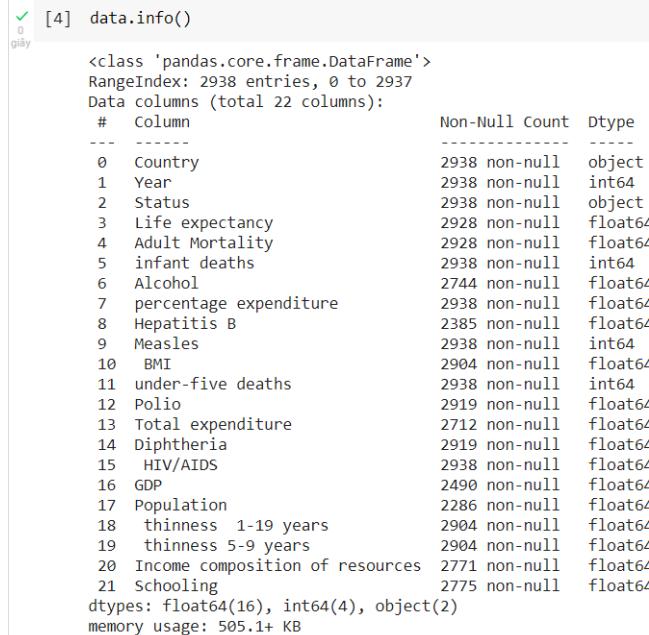
	Country	Year	Status	Life expectancy	Adult Mortality	infant deaths	Alcohol	percentage expenditure	Hepatitis B	Measles	BMI	under-five deaths	Polio	Total expenditure	Diphtheria	HIV/AIDS	GDP	Population
0	Afghanistan	2015	Developing	65.0	263.0	62	0.01	71.279624	65.0	1154	19.1	83	6.0	8.16	65.0	0.1	584.259210	33736494.0
1	Afghanistan	2014	Developing	59.9	271.0	64	0.01	73.523582	62.0	492	18.6	86	58.0	8.18	62.0	0.1	612.696514	327582.0
2	Afghanistan	2013	Developing	59.9	268.0	66	0.01	73.219243	64.0	430	18.1	89	62.0	8.13	64.0	0.1	631.744976	31731688.0
3	Afghanistan	2012	Developing	59.5	272.0	69	0.01	78.184215	67.0	2787	17.6	93	67.0	8.52	67.0	0.1	669.959000	3696958.0
4	Afghanistan	2011	Developing	59.2	275.0	71	0.01	7.097109	68.0	3013	17.2	97	68.0	7.87	68.0	0.1	63.537231	2978599.0
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
2933	Zimbabwe	2004	Developing	44.3	723.0	27	4.36	0.000000	68.0	31	27.1	42	67.0	7.13	65.0	33.6	454.366654	12777511.0
2934	Zimbabwe	2003	Developing	44.5	715.0	26	4.06	0.000000	7.0	998	26.7	41	7.0	6.52	68.0	36.7	453.351155	12633897.0
2935	Zimbabwe	2002	Developing	44.8	73.0	25	4.43	0.000000	73.0	304	26.3	40	73.0	6.53	71.0	39.8	57.348340	125525.0
2936	Zimbabwe	2001	Developing	45.3	686.0	25	1.72	0.000000	76.0	529	25.9	39	76.0	6.16	75.0	42.1	548.587312	12366165.0
2937	Zimbabwe	2000	Developing	46.0	665.0	24	1.68	0.000000	79.0	1483	25.5	39	78.0	7.10	78.0	43.5	547.358879	12222251.0

### 8.2.2 Data exploration

In this section, we will dive deeper to see some insights of this data. This step will give us an overview about which features that we should select and how we can preprocess the data later.

By using `data.info()`, it gives us the name of all columns along with their data type.

**Output:**



```
[4] data.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2938 entries, 0 to 2937
Data columns (total 22 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   Country          2938 non-null    object  
 1   Year              2938 non-null    int64  
 2   Status             2938 non-null    object  
 3   Life expectancy   2928 non-null    float64 
 4   Adult Mortality   2928 non-null    float64 
 5   infant deaths    2938 non-null    int64  
 6   Alcohol            2744 non-null    float64 
 7   percentage expenditure  2938 non-null    float64 
 8   Hepatitis B       2385 non-null    float64 
 9   Measles            2938 non-null    int64  
 10  BMI                2904 non-null    float64 
 11  under-five deaths 2938 non-null    int64  
 12  Polio              2919 non-null    float64 
 13  Total expenditure  2712 non-null    float64 
 14  Diphtheria         2919 non-null    float64 
 15  HIV/AIDS           2938 non-null    float64 
 16  GDP                2490 non-null    float64 
 17  Population          2286 non-null    float64 
 18  thinness 1-19 years 2904 non-null    float64 
 19  thinness 5-9 years  2904 non-null    float64 
 20  Income composition of resources 2771 non-null    float64 
 21  Schooling           2775 non-null    float64 
dtypes: float64(16), int64(4), object(2)
memory usage: 505.1+ KB
```

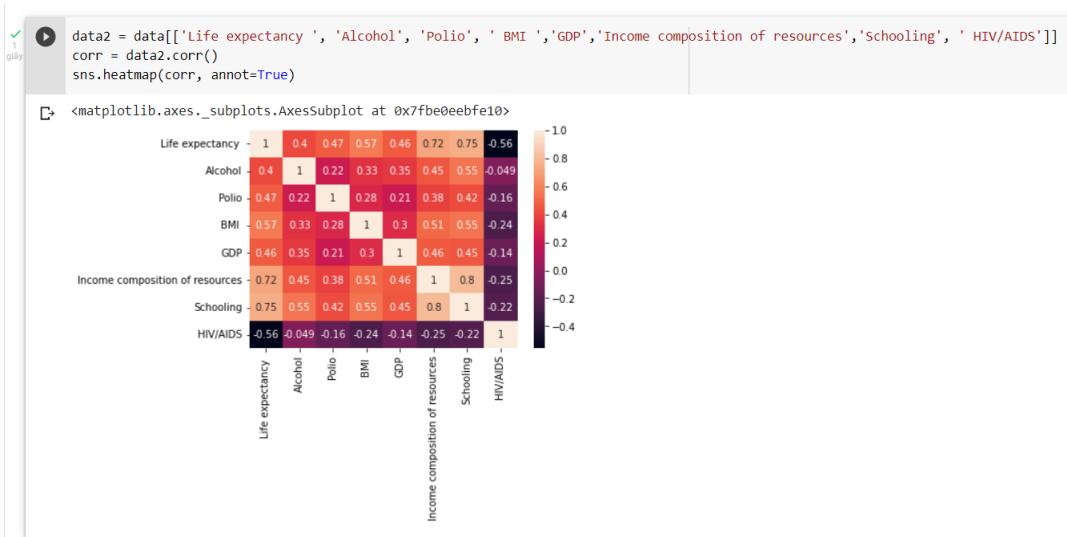
From the output, we can see that there are 22 columns with their own datatypes. However, we also figure out that some null values exists in the data, which we will need to clean later in the data preprocessing part.

For selecting data, we need to know how each feature can correlate linearly to the others. Therefore, we use the heatmap, which will show the Pearson correlation coefficient between every

feature in our dataset.

```
1 data2 = data[['Life expectancy ', 'Alcohol', 'Polio', ' BMI ','GDP','Income
   composition of resources','Schooling', ' HIV/AIDS']]
2 corr = data2.corr()
3 sns.heatmap(corr, annot=True)
```

**Output:**



The correlation coefficient ranges from -1 to 1. An absolute value of exactly 1 implies that a linear equation describes the relationship between X and Y perfectly, with all data points lying on a line. The correlation sign is determined by the regression slope: a value of +1 implies that all data points lie on a line for which Y increases as X increases, and vice versa for -1. A value of 0 implies that there is no linear dependency between the variables.

From the concept of correlation coefficient and the heatmap that we generate above, we can figure out which features we should select. And those are: Status, Alcohol, BMI, Polio, HIV/AIDS, GDP, Income composition of resources and Schooling, which have the highest absolute values of correlation coefficient for Life expectancy.

Now, from the selected features, we create a new table to work on later.

```
1 data3 = data[['Status', 'Life expectancy ', 'Alcohol', 'Polio', ' BMI ','GDP','Income
   composition of resources','Schooling', ' HIV/AIDS']]
2 data3
```

**Output:**



```
[8]: data3 = data[['Status', 'Life expectancy', 'Alcohol', 'Polio', 'BMI', 'GDP', 'Income composition of resources', 'Schooling', 'HIV/AIDS']]  
data3
```

	Status	Life expectancy	Alcohol	Polio	BMI	GDP	Income composition of resources	Schooling	HIV/AIDS
0	Developing	65.0	0.01	6.0	19.1	584.259210		0.479	10.1
1	Developing	59.9	0.01	58.0	18.6	612.696514		0.476	10.0
2	Developing	59.9	0.01	62.0	18.1	631.744976		0.470	9.9
3	Developing	59.5	0.01	67.0	17.6	669.959000		0.463	9.8
4	Developing	59.2	0.01	68.0	17.2	63.537231		0.454	9.5
...	...	...	...	...	...	...	...	...	...
2933	Developing	44.3	4.36	67.0	27.1	454.366654		0.407	9.2
2934	Developing	44.5	4.06	7.0	26.7	453.351155		0.418	9.5
2935	Developing	44.8	4.43	73.0	26.3	57.348340		0.427	10.0
2936	Developing	45.3	1.72	76.0	25.9	548.587312		0.427	9.8
2937	Developing	46.0	1.68	78.0	25.5	547.35879		0.434	9.8

2938 rows × 9 columns

### 8.2.3 Data preprocessing

As said earlier, there are some null values exists in the data. We now need to clean all of it to prepare for our model.

First, let's see how many rows are having null values by using the following code:

```
1 data3.isnull().sum()
```

**Output:**

```
[9]: data3.isnull().sum()  
Status                      0  
Life expectancy              10  
Alcohol                     194  
Polio                       19  
BMI                         34  
GDP                        448  
Income composition of resources 167  
Schooling                   163  
HIV/AIDS                    0  
dtype: int64
```

From the output above, we can see that there are 10 rows having null values for Life expectancy column, 194 for Alcohol, 19 for Polio, 34 for BMI, 448 for GDP, 167 for Income composition of resources and 163 for Schooling.

To clean these rows, we will use the following code, which helps us to remove all of these rows from the table.

```
1 data4 = data3.dropna()  
2 data4.isnull().sum()
```

**Output:**



```
[10] data4 = data3.dropna()  
data4.isnull().sum()  
  
Status 0  
Life expectancy 0  
Alcohol 0  
Polio 0  
BMI 0  
GDP 0  
Income composition of resources 0  
Schooling 0  
HIV/AIDS 0  
dtype: int64
```

Look at the output, we know that all rows in our data are non-null. And let's see how many rows that we still have after cleaning:

The screenshot shows a Jupyter Notebook cell with the code `data4`. Below the code, the resulting DataFrame is displayed. The DataFrame has 2304 rows and 9 columns. The columns are labeled: Status, Life expectancy, Alcohol, Polio, BMI, GDP, Income composition of resources, Schooling, and HIV/AIDS. The data shows various values for each row, such as Life expectancy ranging from 44.3 to 65.0 and Schooling ranging from 9.2 to 10.1.

	Status	Life expectancy	Alcohol	Polio	BMI	GDP	Income composition of resources	Schooling	HIV/AIDS
0	Developing	65.0	0.01	6.0	19.1	584.259210	0.479	10.1	0.1
1	Developing	59.9	0.01	58.0	18.6	612.696514	0.476	10.0	0.1
2	Developing	59.9	0.01	62.0	18.1	631.744976	0.470	9.9	0.1
3	Developing	59.5	0.01	67.0	17.6	669.959000	0.463	9.8	0.1
4	Developing	59.2	0.01	68.0	17.2	63.537231	0.454	9.5	0.1
...	...	...	...	...	...	...	...	...	...
2933	Developing	44.3	4.36	67.0	27.1	454.366654	0.407	9.2	33.6
2934	Developing	44.5	4.06	7.0	26.7	453.351155	0.418	9.5	36.7
2935	Developing	44.8	4.43	73.0	26.3	57.348340	0.427	10.0	39.8
2936	Developing	45.3	1.72	76.0	25.9	548.587312	0.427	9.8	42.1
2937	Developing	46.0	1.68	78.0	25.5	547.358879	0.434	9.8	43.5

The table still have 2304 rows, which is enough for training our model later. But now, there is still one thing we have to do before that. Consider the Status column, we see that it is of object datatype, which we are unable to give the model since the model cannot interpret this datatype. To convert it to numeric type, we use One hot encoding, which is one of the most popular method. Using the `get_dummies` function that the library supports, it becomes very simple:

```
1 dummies = pd.get_dummies(data4.Status)  
2 dummies.head(3)
```

**Output:**

```
[12] dummies = pd.get_dummies(data4.Status)  
dummies.head(3)  
  
Developed  Developed  
0 0 1  
1 0 1  
2 0 1
```



Using this method, our Status column is now divided into two, which are Developed and Developing. These two new columns are also two unique values in the Status column at the beginning. In short, the Status column was converted to 2 columns that contain only numerical values (0 & 1). 0 indicates non-existent while 1 indicates existent.

From here, we just need to concatenate Developed and Developing columns above to our main table. After this step, we will be ready for building our model.

	Life expectancy	Alcohol	Polio	BMI	GDP	Income composition of resources	Schooling	HIV/AIDS	Developed	Developing
0	65.0	0.01	6.0	19.1	584.259210		0.479	10.1	0.1	0
1	59.9	0.01	58.0	18.6	612.696514		0.476	10.0	0.1	0
2	59.9	0.01	62.0	18.1	631.744976		0.470	9.9	0.1	0

## 8.3 Fitting model

### 8.3.1 Data splitting

Before training our model, we need to split our data into training set and testing set. Here, we split them with the ratio 4/5 (20% for testing, 80% for training).

In these two sets, we will assign input features into X and the output (label) to y. To do that, we use the following code:

```
1 #20% test - 80% train
2 X = data5.drop('Life expectancy ', axis='columns')
3 y = data5['Life expectancy ']
4 X_train, X_test, y_train, y_test = train_test_split(X,y,test_size=0.2,random_state
=10)
```

### 8.3.2 Train model

Now, we are able to build a Lasso Regression using the library Scikit-learn with just a few lines of code:

```
1 #Lasso Regression
2 modelLasso = Lasso(alpha=0.8)
3 modelLasso.fit(X_train, y_train)
```

Run the above code and it will train our model with the given data X\_train, y\_train. To see how the model can predict the life expectancy, we have to use our self-defined function instead of the **predict** function that Scikit-learn has provided. Let's call it **predict\_expectancy**:

```
1 def predict_expectancy(Status,Alcohol,Polio,BMI,GDP,Income,Schooling,HIV):
2     loc_index = np.where(X.columns==Status)[0][0]
3     x = np.zeros(len(X.columns))
4     x[0] = Alcohol
5     x[1] = Polio
6     x[2] = BMI
7     x[3] = GDP
8     x[4] = Income
9     x[5] = Schooling
10    x[6] = HIV
11    if loc_index >= 0:
```



```
12     x[loc_index] = 1
13     return modelLasso.predict([x])[0]
14
15 predict_expectancy('Developed', 0.01, 6.0, 19.1, 584.259210, 0.479, 10.1, 0.1)
```

We gain the output:

```
[✓] [77] predict_expectancy('Developed', 0.01, 6.0, 19.1, 584.259210, 0.479, 10.1, 0.1)
0 giây
61.81712205829855
```

### 8.3.3 Evaluation

After building the Lasso model, we now use some metrics to evaluate how good our model is. Then, we can compare our Lasso Regression model with some other models to see the difference. First, let's just evaluate the Lasso model. We are going to use K-fold cross validation to this:

```
1 cv = ShuffleSplit(n_splits = 5, test_size = 0.2, random_state = 0)
2 cross_val_score(Lasso(alpha=0.8), X, y, cv = cv)
```

Output:

```
[✓] [70] cv = ShuffleSplit(n_splits = 5, test_size = 0.2, random_state = 0)
0 giây
cross_val_score(Lasso(alpha=0.8), X, y, cv = cv)

array([0.77560974, 0.7813822 , 0.77096243, 0.76938786, 0.74200145])
```

Also, let's consider some other metrics such as MAE, MSE, RMSE:

Output:

```
[✓] [70] #Lasso Regression evaluation
0 giây
test_pred = modelLasso.predict(X_test)
train_pred = modelLasso.predict(X_train)
print("Lasso Regression model")
print('Test set evaluation:\n_____')
print_evaluate(y_test, test_pred)
print('Train set evaluation:\n_____')
print_evaluate(y_train, train_pred)

↳ Lasso Regression model
Test set evaluation:
_____
MAE: 3.4437549579484124
MSE: 20.113241025529117
RMSE: 4.484778815675208
R2 Square 0.7489349293449861

Train set evaluation:
_____
MAE: 3.4626282149783387
MSE: 20.556543305350242
RMSE: 4.53393243281704
R2 Square 0.7897641448891335
```



From the output above, we can see that our model gets the score of around 77%, which are quite high. But compared to some other models, it may be a little bit lower. However, this does not really matter, let's see why in the following section.

#### 8.3.4 Comparison

In this section, we will use the same metrics that have applied for Lasso model as above for two following models: Linear Regression and Random Forest Regressor. Let's see the result.

Random Forest Regressor evaluation:

```
✓ [22] #Random Forest Regressor evaluation
0
giây
    test_pred = rf_reg.predict(X_test)
    train_pred = rf_reg.predict(X_train)
    print('Test set evaluation:\n_____')
    print_evaluate(y_test, test_pred)
    print('Train set evaluation:\n_____')
    print_evaluate(y_train, train_pred)

Test set evaluation:
_____
MAE: 1.4096761388286234
MSE: 4.3045077609327835
RMSE: 2.0747307683005003
R2 Square 0.9462686523886463

Train set evaluation:
_____
MAE: 0.5833295713510787
MSE: 0.7569218832989897
RMSE: 0.8700125765177131
R2 Square 0.9922588094202553
```

Linear Regression evaluation:

```
✓ [23] #Linear Regression evaluation
0
giây
    test_pred = lin_reg.predict(X_test)
    train_pred = lin_reg.predict(X_train)
    print('Test set evaluation:\n_____')
    print_evaluate(y_test, test_pred)
    print('Train set evaluation:\n_____')
    print_evaluate(y_train, train_pred)

Test set evaluation:
_____
MAE: 3.194931421790469
MSE: 18.400971697366654
RMSE: 4.289635380468444
R2 Square 0.7703084622982219

Train set evaluation:
_____
MAE: 3.2628688786258406
MSE: 18.724129432845753
RMSE: 4.327138711995
R2 Square 0.8085046058548022
```

From the above results, the Linear Regression model has a little better score than the Lasso model, while the Random Regression has the highest score of over 90%. However, as we have said before about Lasso Regression model, the core of this model is that it improves the issue of



overfitting and give the estimation as closed as possible to the real life cases. Therefore, a lower score does not mean that it is worse than two other models.

## 9 Medical Cost Personal - Mạnh Hùng

### 9.1 Background and Data Overview

One of the most important sectors of the global economy is health care. Health-care spending amounted for 9.95 percent of global gross domestic product in 2014, according to the World Bank (GDP). In addition, over the previous ten years, per capita health expenditures have risen. The Centers for Medicare Medicaid Services (CMS) in the United States reported that health care accounted for 17.5 percent of national GDP in 2014. Over the next few years, this figure is likely to rise.

A health insurance business can only make money if it collects more money than it spends on its beneficiaries' medical treatment. Medical costs, on the other hand, are difficult to forecast because the majority of money comes from patients with unusual diseases. The goal of this article is to accurately anticipate insurance prices based on individuals' data, such as age, BMI, whether or not they smoke, and so on. These estimations could be used to develop actuarial tables that adjust the price of annual premiums based on predicted treatment costs. This is a case of regression.

We are going to take a overview about this data set. For predicting health insurance costs, we utilize [Miri Choi's Medical Cost Personal Datasets](#) hosted on Kaggle. This is a good dataset for regression models in general or LASSO in particular. The column descriptions look like this:

- **age:** age of primary beneficiary
- **sex:** insurance contractor gender, female, male
- **bmi:** Body mass index, providing an understanding of body, weights that are relatively high or low relative to height, objective index of body weight ( $\text{kg} / \text{m}^2$ ) using the ratio of height to weight, ideally 18.5 to 24.9
- **children:** Number of children covered by health insurance / Number of dependents
- **smoker:** Yes/No
- **region:** the beneficiary's residential area in the US, northeast, southeast, southwest, northwest.
- **charges:** Individual medical costs billed by health insurance

### 9.2 Implementation

#### 9.2.1 Import package and acquire the data

Firstly, we need to import some libraries that we are going to work with. All of the relevant libraries are listed below:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import pandas as pd
4 import seaborn as sns
5 from sklearn import metrics
6 from sklearn.model_selection import GridSearchCV
7 from sklearn.linear_model import Lasso
8 from sklearn.tree import DecisionTreeRegressor
```



```
9 from sklearn.linear_model import LinearRegression
10 from sklearn.model_selection import ShuffleSplit
11 from sklearn.model_selection import cross_val_score
12 from sklearn.linear_model import SGDRegressor
```

Once we download the CSV data, we can import it using read\_csv.

```
1 %%shell
2 # RUN THIS CELL
3
4 ## Guidelines for sharing files
5 ## 1. share a file (here, zip) with "anyone with link" in GoogleDrive
6 ## 2. copy the link, here: https://drive.google.com/file/d/1EHMGaGbf-
7 ## GGNevLqaRNH13FbupJQOECN/view?usp=sharing
8 ## 3. extract the file-id, here: 1EHMGaGbf-GGNevLqaRNH13FbupJQOECN
9 ## 4. download the file-id with gdown, as follows
10 IMAGES="1EHMGaGbf-GGNevLqaRNH13FbupJQOECN"
11 cd /content/
12 gdown -q --id $IMAGES
```

We then use head() to sample the data.

```
1 df = pd.read_csv("insurance.csv")
2 df.head()
```

**Output:**

```
   ↗    age      sex      bmi  children  smoker      region      charges
0     19  female  27.900          0     yes  southwest  16884.92400
1     18    male  33.770          1     no  southeast  1725.55230
2     28    male  33.000          3     no  southeast  4449.46200
3     33    male  22.705          0     no northwest  21984.47061
4     32    male  28.880          0     no northwest  3866.85520
      ↗    age      sex      bmi  children  smoker      region      charges
1333    50    male  30.97          3     no northwest  10600.5483
1334    18  female  31.92          0     no northeast  2205.9808
1335    18  female  36.85          0     no southeast  1629.8335
1336    21  female  25.80          0     no southwest  2007.9450
1337    61  female  29.07          0     yes northwest  29141.3603
```

### 9.2.2 Preprocessing data

In this step we will drop null values if any, then we will visualize our dataset by plotting correlation. Next we select features based on that correlation and convert object-type data into a number using one-hot encoding technique.

#### 9.2.2.a Check Null values

```
1 Total_null_values = df.isnull().sum()
2 print(Total_null_values)
```

**Output:**



```
↳ age      0
    sex     0
    bmi     0
    children  0
    smoker   0
    region   0
    charges  0
dtype: int64
```

As we can see, there is no null value. So we do not have to drop any row. Instead, i use `data.info()`, it gives us the name of all columns along with their data type.

```
1 df.shape
2 df.info()
```

#### Output:

```
↳ <class 'pandas.core.frame.DataFrame'>
RangeIndex: 1338 entries, 0 to 1337
Data columns (total 7 columns):
 #   Column   Non-Null Count  Dtype  
--- 
 0   age       1338 non-null   int64  
 1   sex       1338 non-null   object  
 2   bmi       1338 non-null   float64 
 3   children  1338 non-null   int64  
 4   smoker    1338 non-null   object  
 5   region    1338 non-null   object  
 6   charges   1338 non-null   float64 
dtypes: float64(2), int64(2), object(3)
memory usage: 73.3+ KB
```

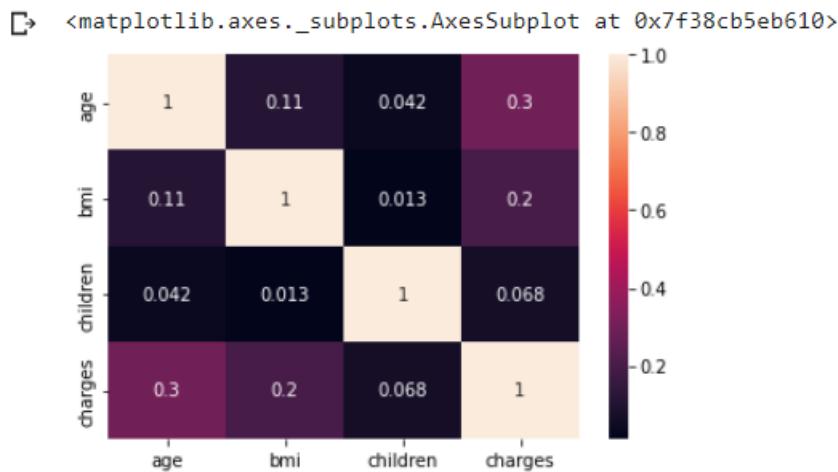
From the picture above, there are 7 columns with their own datatypes and 1338 rows.

#### 9.2.2.b Visualization

For selecting data, we need to know how each feature can correlate linearly to the others. Therefore, we use the heatmap, which will show the Pearson correlation coefficient between every feature in our dataset:

```
1 corr = df.corr()
2 sns.heatmap(corr, annot=True)
```

#### Output:



The correlation coefficient ranges from -1 to 1. An absolute value of exactly 1 implies that a linear equation describes the relationship between X and Y perfectly, with all data points lying on a line. The correlation sign is determined by the regression slope: a value of +1 implies that all data points lie on a line for which Y increases as X increases, and vice versa for -1. A value of 0 implies that there is no linear dependency between the variables.

From the concept of correlation coefficient and the heatmap above, we can figure out which features we should select. And those are: "age", "bmi", which have the highest absolute values of correlation coefficient for "charges". And we also take "sex", "smoker", "region" features too.

### 9.2.2.c Select features (based on correlation)

Now, from the selected features, We make a new data frame on which we will work later.

```
1 df2 = df[['age', 'sex', 'bmi', 'smoker', 'region', 'charges']]  
2 df2
```

**Output:**

	age	sex	bmi	smoker	region	charges
0	19	female	27.900	yes	southwest	16884.92400
1	18	male	33.770	no	southeast	1725.55230
2	28	male	33.000	no	southeast	4449.46200
3	33	male	22.705	no	northwest	21984.47061
4	32	male	28.880	no	northwest	3866.85520
...	...	...	...	...	...	...
1333	50	male	30.970	no	northwest	10600.54830
1334	18	female	31.920	no	northeast	2205.98080
1335	18	female	36.850	no	southeast	1629.83350
1336	21	female	25.800	no	southwest	2007.94500
1337	61	female	29.070	yes	northwest	29141.36030

1338 rows × 6 columns

From the table above, we can see that the data is made up of two forms: Numerical and Categorical.

Building a model with categorical data is hard but not impossible. For simplicity, we proceed to convert categorical data into numerical data. For this purpose, we use the **One hot encoding** technique.

One hot encoding is a technique where we replace the categorical data with binary digits. The categorical column is split into the same number of columns as the values. The respective column is then given a '1' or a '0' corresponding to the values.

We use one-hot encoding by using `get_dummies()`:

```
1 dummies_sex = pd.get_dummies(df2.sex)
2 print(dummies_sex.head(2))
3 dummies_smoker = pd.get_dummies(df2.smoker)
4 print(dummies_smoker.head(2))
5 dummies_region = pd.get_dummies(df2.region)
6 print(dummies_region.head(4))
```

**Output:**

```
↳   female  male
  0      1    0
  1      0    1
    no yes
  0    0    1
  1    1    0
    northeast northwest southeast southwest
  0      0      0      0      1
  1      0      0      1      0
  2      0      0      1      0
  3      0      1      0      0
```

Next, we create a new data frame containing our numerical data only.

```
1 df3 = pd.concat([df2.drop(['sex', 'smoker', 'region'], axis='columns'), dummies_sex
                  , dummies_smoker, dummies_region], axis='columns')
2 df3.head(3)
```

**Output:**

```
↳   age  bmi  charges  female  male  no  yes  northeast  northwest  southeast  southwest
  0  19  27.90  16884.9240      1    0    0    1      0      0      0      1
  1  18  33.77  1725.5523      0    1    1    0      0      0      1      0
  2  28  33.00  4449.4620      0    1    1    0      0      0      1      0
```

## 9.3 Building the model and evaluation

### 9.3.1 Data splitting

We'll assign input features to X and output (label) to y in these two sets. We use the following code to accomplish this:

```
1 X = df3.drop('charges', axis='columns')
2 X.head()
```

**Output:**



	age	bmi	female	male	no	yes	northeast	northwest	southeast	southwest
0	19	27.900	1	0	0	1	0	0	0	1
1	18	33.770	0	1	1	0	0	0	1	0
2	28	33.000	0	1	1	0	0	0	1	0
3	33	22.705	0	1	1	0	0	1	0	0
4	32	28.880	0	1	1	0	0	1	0	0

```
1 y = df3['charges']
2 y.head()
```

#### Output:

```
0    16884.92400
1    1725.55230
2    4449.46200
3    21984.47061
4    3866.85520
Name: charges, dtype: float64
```

We must divide our data into training and testing sets before training our model. We split them in this case in a 4/5 ratio (20% for testing, 80% for training).

```
1 #20% test - 80% train
2 from sklearn.model_selection import train_test_split
3 X_train, X_test, y_train, y_test = train_test_split(X,y,test_size=0.2,random_state
=10)
```

#### 9.3.2 Train model

Now, we are able to build a Lasso Regression using the library Scikit-learn with just a few lines of code:

```
1 from sklearn.linear_model import Lasso
2 modelLasso = Lasso(alpha=0.8)
3 modelLasso.fit(X_train, y_train)
```

Now our model is ready to use.

I also create a function to calculate the predict output, which is the money spent given the input.

```
#age  bmi  smoker(no yes)
def predict_medical_cost(x_test):
    print(len(X.columns))
    x = np.zeros(len(X.columns))
    x[0] = x_test[0] #age
    x[1] = x_test[1] #bmi
    x[2] = x_test[2] == 'no' #smoker == 'no'
    x[3] = x_test[2] == 'yes' #smoker == 'yes'
    print(x)
    return modelLasso.predict([x])[0]
x_to_test = (28, 33.000, "no")
y_predict = predict_medical_cost(x_to_test)
print(y_predict)

4
[28. 33. 1. 0.]
6168.7624600678355
```

## 9.4 Evaluation and Comparision

### 9.4.1 Evaluation

After we've built the Lasso model, we'll utilize some metrics to see how good it is. We can then compare our Lasso Regression model to a few other models to see how they differ.

First of all, we will focus on evaluating the Lasso model using K-fold cross validation:

```
[32] cv = ShuffleSplit(n_splits= 5, test_size= 0.2, random_state= 0)
      cross_val_score(Lasso(alpha = 1), X, y, cv = cv)

array([0.79453718, 0.74563532, 0.70688027, 0.77266459, 0.81012932])
```

### 9.4.2 Comparision

Before comparing with other methods like linear regression, decision tree, ... We will check the score of test set and train set by using this function:



```
[36] def print_evaluate(true, predicted):
    mae = metrics.mean_absolute_error(true, predicted)
    mse = metrics.mean_squared_error(true, predicted)
    rmse = np.sqrt(metrics.mean_squared_error(true, predicted))
    r2_square = metrics.r2_score(true, predicted)
    print('MAE:', mae)
    print('MSE:', mse)
    print('RMSE:', rmse)
    print('R2 Square', r2_square)
    print('_____')

def evaluate(true, predicted):
    mae = metrics.mean_absolute_error(true, predicted)
    mse = metrics.mean_squared_error(true, predicted)
    rmse = np.sqrt(metrics.mean_squared_error(true, predicted))
    r2_square = metrics.r2_score(true, predicted)
    return mae, mse, rmse, r2_square
```

And then, I compare score from R squared, MSE and MAE:

```
[ ] test_pred = modelLasso.predict(X_test)
train_pred = modelLasso.predict(X_train)

print('Test set evaluation:\n_____')
print_evaluate(y_test, test_pred)
print('Train set evaluation:\n_____')
print_evaluate(y_train, train_pred)
```

Test set evaluation:

---

```
MAE: 4575.575844024035
MSE: 42876310.32570865
RMSE: 6548.000483025994
R2 Square 0.6942880783720784
```

---

Train set evaluation:

---

```
MAE: 4136.597783092333
MSE: 35559732.06997611
RMSE: 5963.198141096445
R2 Square 0.7598838051825363
```

---

Next, we compare the Lasso with the other two models, which are decision tree and Linear Regression. We can see that the score of lasso is quite good compared to the two algorithms.



	model	best_score	best_params
0	linear_regression	0.765973	{'normalize': False}
1	lasso	0.765971	{'alpha': 0.5, 'selection': 'random'}
2	decision_tree	0.711208	{'criterion': 'friedman_mse', 'splitter': 'ran...}



## 10 Predict final grade - Chương

### 10.1 Data Overview

This data approach student achievement in secondary education of two Portuguese schools. The data attributes include student grades, demographic, social and school related features) and it was collected by using school reports and questionnaires. Attribute information:

- **sex:** student's sex ('F': female, 'M': male)
- **age:** student's age (numeric: from 15 to 22)
- **studytime:** weekly study time (numeric: 1 - <2 hours, 2 - 2 to 5 hours, 3 - 5 to 10 hours, or 4- >10 hours)
- **failures:** number of past class failures (numeric: n if  $1 \leq n \leq 3$ , else 4)
- **higher:** wants to take higher education (binary: yes or no)
- **absences:** number of school absences (numeric: from 0 to 93)
- **charges:** Individual medical costs billed by health insurance
- **G1:** first period grade (numeric: from 0 to 20)
- **G2:** second period grade (numeric: from 0 to 20)
- **G3:** final grade (numeric: from 0 to 20, output target)

### 10.2 Import libraries and preprocessing data

#### 10.2.1 Set up libraries and load data

To work with the dataset from csv file and display the summary of the dataset, I'll add some libraries to do that. Firstly, I install the pandas-profiling from github by pip in python:

```
1 ! pip install https://github.com/pandas-profiling/pandas-profiling/archive/master.zip
```

After that, I import those libraries into the google colaboratory file:

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 import numpy as np
5 import pylab
6 import scipy.stats as stats
7 from pandas_profiling import ProfileReport
8
9 from sklearn.metrics import r2_score, mean_absolute_error, mean_squared_error,
   accuracy_score, classification_report, confusion_matrix,
   precision_recall_fscore_support
10 from sklearn.model_selection import train_test_split, cross_val_score
11 from sklearn.naive_bayes import MultinomialNB
12 from sklearn.preprocessing import StandardScaler, PolynomialFeatures
13 from sklearn.linear_model import LinearRegression, Lasso, LassoCV
14 from sklearn.mixture import GaussianMixture
```



```
15 from mlxtend.preprocessing import TransactionEncoder
16 from pandas_profiling import ProfileReport
17 from sklearn.ensemble import AdaBoostClassifier, RandomForestClassifier,
    RandomForestRegressor
18 from sklearn.tree import DecisionTreeClassifier, DecisionTreeRegressor
19 from sklearn.linear_model import LogisticRegression
20 from sklearn.naive_bayes import BernoulliNB, GaussianNB
21 from scipy.stats.stats import pearsonr
22 from sklearn.feature_selection import SelectKBest, f_regression
23 from sklearn import svm
24 from scipy.integrate import odeint
25 import xgboost as xgb
26
27 import os, sys, glob, datetime, warnings
28 from google.colab import files, drive
29 warnings.filterwarnings("ignore")
30 %matplotlib inline
```

After mounting google drive or using !gdown to download dataset, I'll read the dataset from file .csv, I use Pandas to do this task, then I will show some information about dataset by head(), info(), describe() commands, or i can see it through pandas profiling by typing command "profile = ProfileReport(data)", then pandas profiling will generate an html page by itself so i can see all basic information related to data.

```
1 data = pd.read_csv("grade.csv")
2 data.head()
```

```
1 data.describe()
```

```
✓ [38] data = pd.read_csv("grade.csv")
In [38]: data.head(5)

Unamed: 0 school sex age address famsize Pstatus Medu Fedu Mjob Fjob reason guardian traveltim studytime failures schoolsup famsup paid activities nursery higher internet romant
0 1 GP F 18 U GT3 A 4 4 at_home teacher course mother 2 2 0 yes no no no no yes yes no
1 2 GP F 17 U GT3 T 1 1 at_home other course father 1 2 0 no yes no no no yes yes yes
2 3 GP F 15 U LE3 T 1 1 at_home other other mother 1 2 3 yes no yes no yes yes yes yes
3 4 GP F 15 U GT3 T 4 2 health services home mother 1 3 0 no yes yes yes yes yes yes yes
4 5 GP F 16 U GT3 T 3 3 other other home father 1 2 0 no yes yes no yes yes no
```

```
✓ [51] data.describe()

Unamed: 0 age Medu Fedu traveltim studytime failures famrel freetime gout Dalc Walc health absences G1 G2 G3
count 395.000000 395.000000 395.000000 395.000000 395.000000 395.000000 395.000000 395.000000 395.000000 395.000000 395.000000 395.000000 395.000000 395.000000 395.000000 395.000000
mean 198.000000 16.696203 2.749367 2.521519 1.448101 2.035443 0.334177 3.944304 3.235443 3.108861 1.481013 2.291139 3.554430 5.708861 10.908861 10.717949 10.415190
std 114.170924 1.276043 1.994735 1.088201 0.697505 0.839240 0.743651 0.896659 0.998862 1.113278 0.890741 1.287897 1.390303 8.003096 3.319195 3.737868 4.581443
min 1.000000 15.000000 0.000000 0.000000 1.000000 1.000000 0.000000 1.000000 1.000000 1.000000 1.000000 1.000000 1.000000 0.000000 3.000000 0.000000 0.000000
25% 99.500000 16.000000 2.000000 2.000000 1.000000 1.000000 0.000000 4.000000 3.000000 2.000000 1.000000 1.000000 1.000000 4.000000 8.000000 9.000000 8.000000
50% 198.000000 17.000000 3.000000 2.000000 1.000000 2.000000 0.000000 4.000000 3.000000 3.000000 1.000000 2.000000 4.000000 4.000000 11.000000 11.000000 11.000000
75% 296.500000 18.000000 4.000000 3.000000 2.000000 2.000000 0.000000 5.000000 4.000000 4.000000 2.000000 3.000000 5.000000 8.000000 13.000000 13.000000 14.000000
max 395.000000 22.000000 4.000000 4.000000 4.000000 3.000000 5.000000 5.000000 5.000000 5.000000 5.000000 5.000000 75.000000 19.000000 19.000000 20.000000 20.000000
```

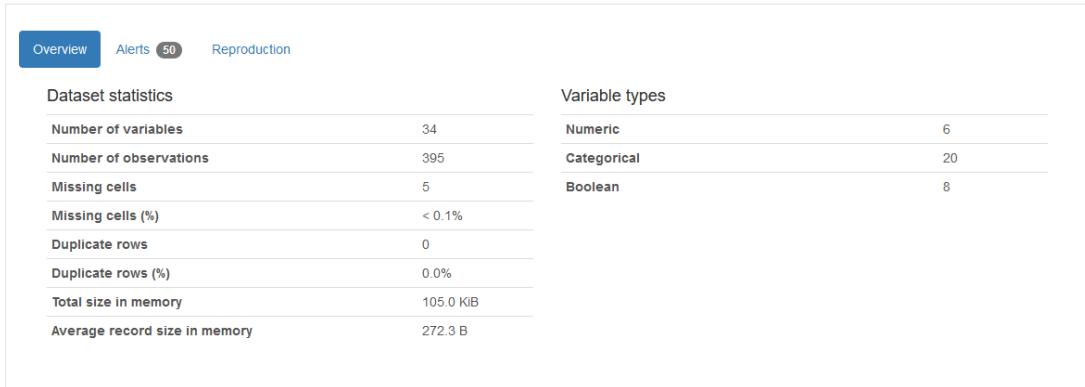
### 10.2.2 Preprocessing

From pandas-profiling, I can see everything about the dataset including null values, or I can use the command .isna().sum() to determine the null value in any columns.

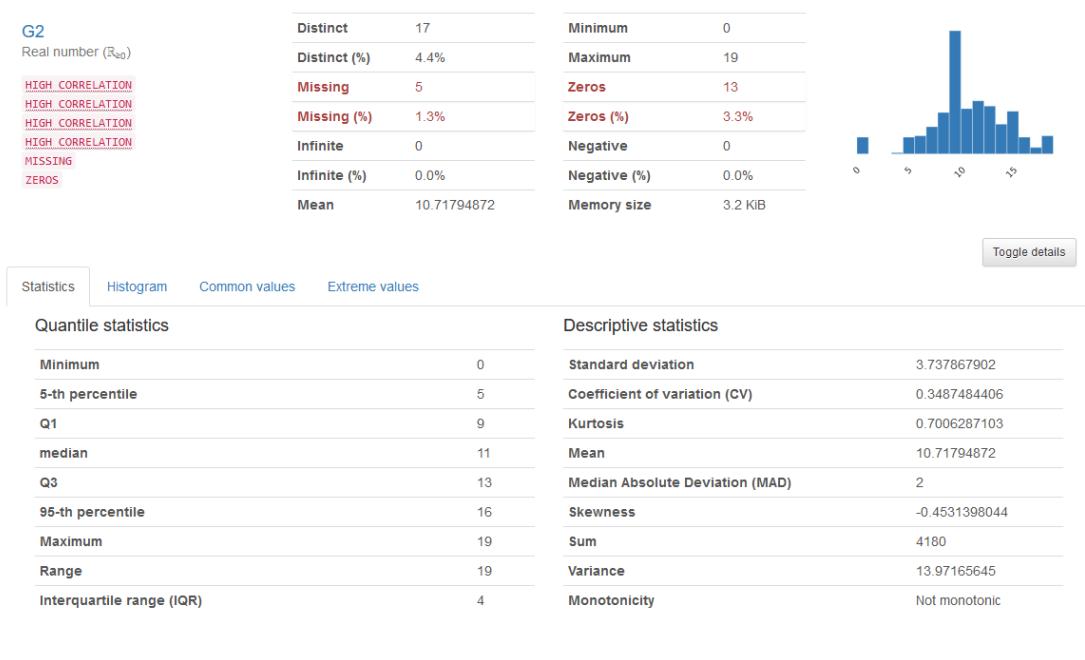
```
1 profile = ProfileReport(data)
2 profile
```



## Overview



I see that in column G2, there are 5 null values, so I need to drop these null values' rows to avoid it may change my prediction results, another way can handle this is we can replace these null values with median or mean.



```
1 data["G2"].fillna(int(data["G2"].median()), inplace = True)
```

Re-check null value in G2 column:

```
1 data["G2"].isnull().sum()
```

And output is 0.

```
1 0
```

Moving on, I noticed that there are some columns that take Yes/No values or string values, to include in the predictive model I need to convert these string values to numeric values (0 or 1)

to get can be calculated.

```

1 data['sex'] = data['sex'].map({'F': 1, 'M': 0})
2 data['schoolsup'] = data['schoolsup'].map({'yes': 1, 'no': 0})
3 data['famsup'] = data['famsup'].map({'yes': 1, 'no': 0})
4 data['paid'] = data['paid'].map({'yes': 1, 'no': 0})
5 data['activities'] = data['activities'].map({'yes': 1, 'no': 0})
6 data['nursery'] = data['nursery'].map({'yes': 1, 'no': 0})
7 data['higher'] = data['higher'].map({'yes': 1, 'no': 0})
8 data['internet'] = data['internet'].map({'yes': 1, 'no': 0})
9 data['romantic'] = data['romantic'].map({'yes': 1, 'no': 0})

```

## 10.3 Visualization and selecting attributes

### 10.3.1 Visualization

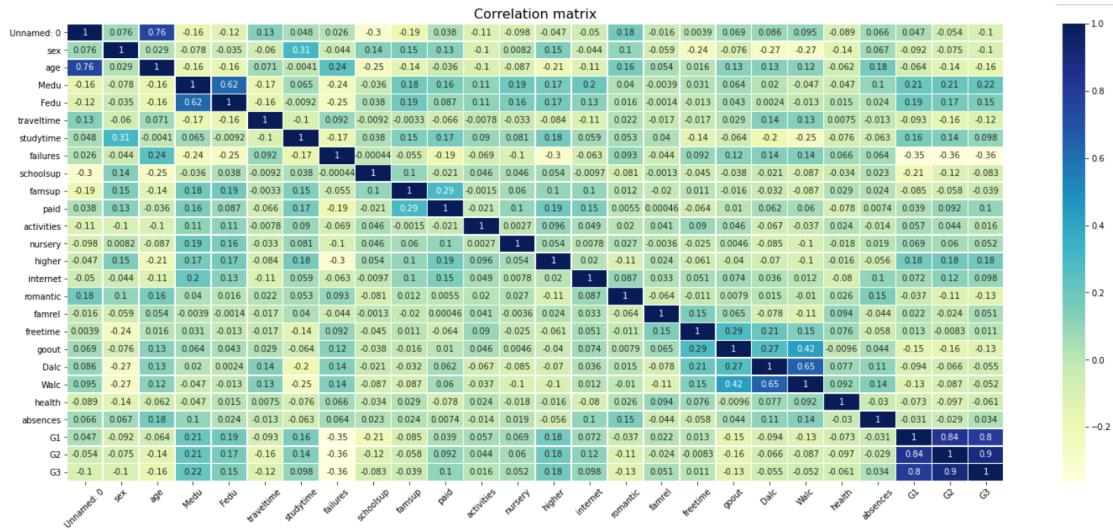
I will plot the correlation matrix to determine that which attributes mostly affect my prediction column (G3). This task is very easy because Pandas-profiling already has support and automatically plots correlation matrices with various types (Spearman's  $\rho$ , Pearson's  $r$ , Kendall's  $\tau$ , Cramér's  $V$  ( $\varphi_c$ ), Phik ( $\varphi_k$ )).





Or I can use the following code to plot the correlation matrix (Pearson) and get the same result

```
1 corr_matrix = data.corr()  
2  
3 fig, ax = plt.subplots(figsize=(25,10))  
4 sns.heatmap(corr_matrix, cmap="YlGnBu", linewidths=.5, annot=True)  
5 plt.title("Correlation matrix", fontsize=16)  
6 plt.xticks(rotation=45)  
7 plt.show()
```

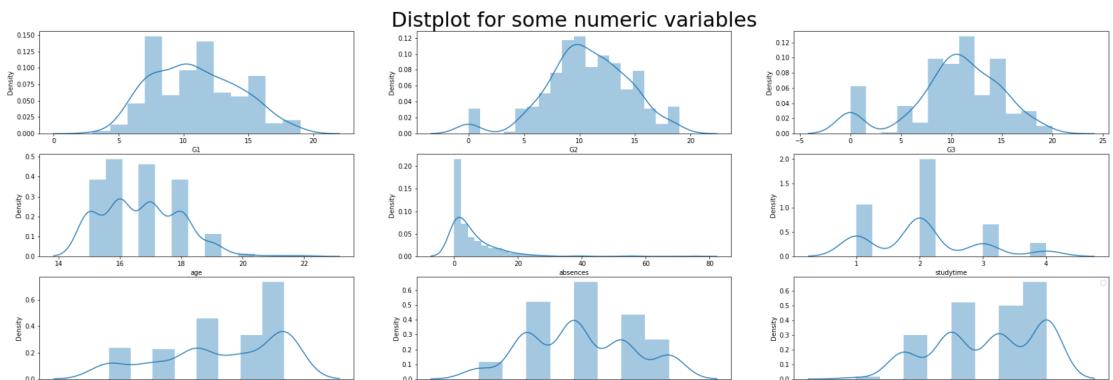


Moreover, I can plot distribution for each attributes by these lines:

```

1 fig, axs = plt.subplots(ncols=3, nrows=3, figsize = (30, 10))
2 sns.distplot(data["G1"], ax=axs[0][0])
3 sns.distplot(data["G2"], ax=axs[0][1]).set_title('Distplot for some numeric
4     variables', fontsize = 32)
5 sns.distplot(data["G3"],ax=axs[0][2])
6 sns.distplot(data["age"], ax=axs[1][0])
7 sns.distplot(data["absences"], ax=axs[1][1])
8 sns.distplot(data["studytime"], ax=axs[1][2])
9 sns.distplot(data["health"], ax=axs[2][0])
10 sns.distplot(data["goout"], ax=axs[2][1])
11 sns.distplot(data["Medu"], ax=axs[2][2])
12 plt.legend()
13 plt.show()

```



### 10.3.2 Selecting attributes and split train/test

From the correlation matrix above, I can see that these attributes (G1, G2, medu, higher) effect mostly to G3, so I choose these attributes to be my input in my model.

```
X = data[['G1',
```



```
2     'G2',
3     'Medu',
4     'higher']] #input
5 y = data[['G3']] #output
```

To split training and testing dataset, I'll use the `train_test_split` to split the data to 75 percent for training and 25 percent for testing.

```
1 X_train, X_test, y_train, y_test = train_test_split(X.values, y.values, test_size
=0.25 ,random_state = 42)
```

## 10.4 Build and evaluate models

### 10.4.1 Lasso Regression

I will use Lasso with cross validation (`LassoCV`) provided by `sklearn` to predict the final grade with the parameters looking like this:

```
1 alphas = np.logspace(-3, 10, 10000)
2 max_iter = 10000
3 model = LassoCV(alphas = alphas, cv = 10, normalize=True, max_iter = max_iter)
4 model.fit(X_train, y_train)
```

and I'll print the intercept and coefficient of the model:

```
1 print(model.intercept_)
2 print(model.coef_)
```

outputs are:

```
1 -2.0475621075482096
2 [0.23318683 0.90733081 0. 0.19104191]
```

To evaluate the model, I use R squared, MSE/RMSE or MAE:

```
1 print("Variance score: ", model.score(X,y)) #R^2
2 print("The Train Score: ",model.score(X_train, y_train))
3 print("The Test Score: ",model.score(X_test, y_test))
4 print("MSE (on train): ", mean_squared_error(y_train, model.predict(X_train)))
5 print("MSE: ", mean_squared_error(y_pred, y_test))
6 print("RMSE: ", np.sqrt(mean_squared_error(y_pred, y_test)))
7 print("MAE: ", mean_absolute_error(y_pred, y_test))
```

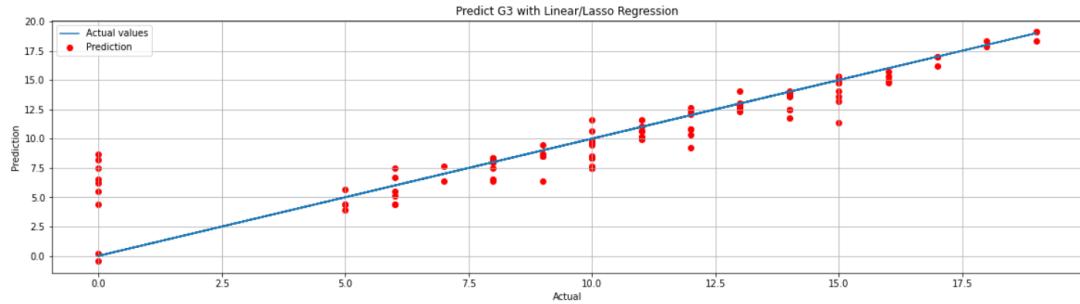
The results are:

```
1 Variance score:  0.8105855104136546
2 The Train Score:  0.8106303290798009
3 The Test Score:  0.8103792428765757
4 MSE (on train):  3.8309577797885175
5 MSE:  4.368455181504587
6 RMSE:  2.090084969924569
7 MAE:  1.250871763911041
```

We can see that the difference R squared on the train and test sets is very small, and the variance score is not too big of a difference and the ratio is quite high at about 81 percent, so it can be concluded that the Lasso model is suitable for prediction on this dataset.

```
1 plt.figure(figsize = (20, 5))
2 plt.plot(y_test, y_test, label="Actual values")
3 plt.scatter(y_test,y_pred, color="red", label="Prediction")
```

```
4 plt.xlabel("Actual")
5 plt.ylabel("Prediction")
6 plt.title("Predict G3 with Linear/Lasso Regression")
7 plt.grid()
8 plt.legend()
9 plt.show()
```



#### 10.4.2 Random Forest and XGB

I'll do the same thing with Random Forest and XGB to compare 3 models later:

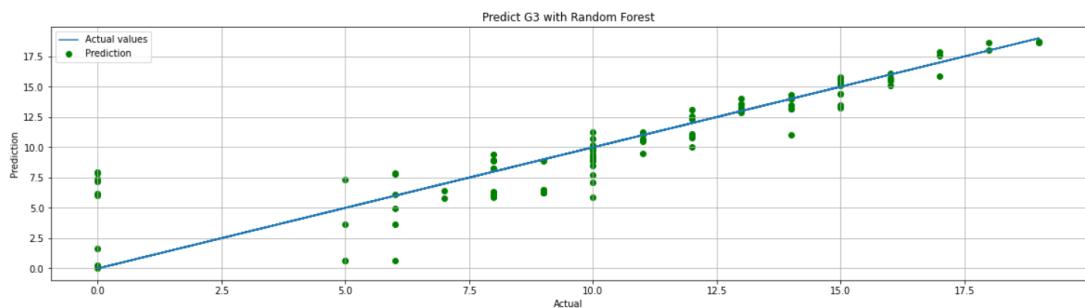
```
1 rrf = RandomForestRegressor(n_estimators=100)
2 model_rrf = rrf.fit(X_train, y_train)
3 y_rrf_pred = model_rrf.predict(X_test)

1 xgb_model = xgb.XGBRegressor(random_state=42)
2 xgb_model.fit(X_train, y_train)
3 y_pred_xgb = xgb_model.predict(X_test)
```

Here is the score from Random Forest model:

```
1 Variance score:  0.8657898120410552
2 The Train Score:  0.903109795283467
3 The Test Score:  0.767743550156823
4 MSE (on train):  1.9600936186899816
5 MSE:  5.350690014885247
6 RMSE:  2.3131558561595558
7 MAE:  1.4724940676607343
```

We can see that Random Forest tend to get the overfitting because the score on train is higher than variance score and test score, so when we use this model, prediction on the training set is good but it's not suitable for the real prediction.



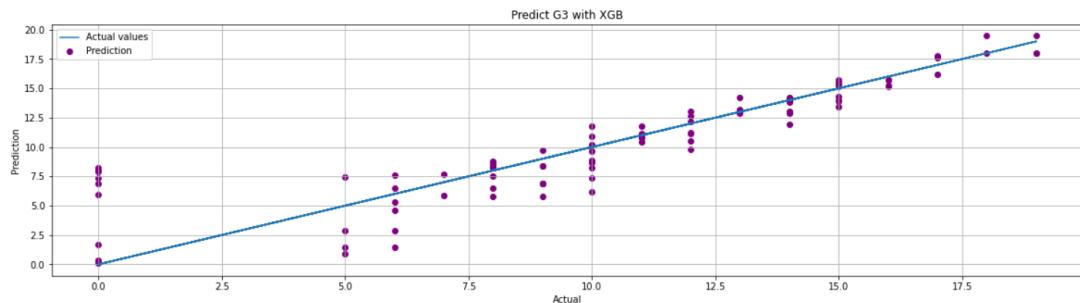
And here is the score from XGB model:

```

1 Variance Score:  0.8420719906587261
2 The Train Score:  0.8675956554429577
3 The Test Score:  0.7749858283072458
4 MSE (on train):  2.678546418725897
5 MSE:  5.183843473440847
6 RMSE:  2.276805541420006
7 MAE:  1.4065777936367074

```

It's overfitting too.



## 10.5 Comparison and new prediction

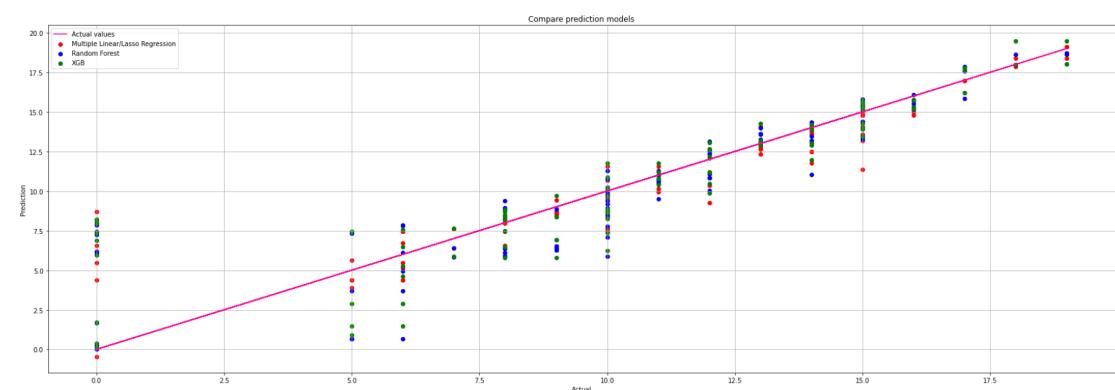
### 10.5.1 Compare 3 models

Firstly, I'll plot the prediction value of each models and compare it to the actual values.

```

1 plt.figure(figsize=(30,10))
2 plt.plot(y_test, y_test, label="Actual values", color="#fc0398")
3 plt.scatter(y_test, y_pred, color='red', label="Multiple Linear/Lasso Regression")
4 plt.scatter(y_test, df_rff['Prediction'], color='blue', label="Random Forest")
5 plt.scatter(y_test, df_xgb['Prediction'], color='green', label="XGB")
6 plt.xlabel("Actual")
7 plt.ylabel("Prediction")
8 plt.title("Compare prediction models")
9 plt.grid()
10 plt.legend()
11 plt.show()

```



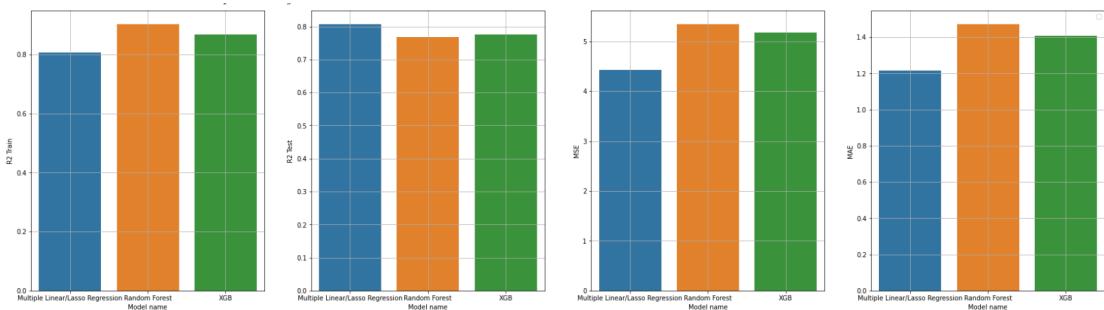
And then, I compare score from R squared, MSE and MAE:

```
1 models = ["Multiple Linear/Lasso Regression", "Random Forest", "XGB"]
```

```
2 errors_mse = [mean_squared_error(y_pred, y_test), mean_squared_error(y_rrf_pred,
3   y_test), mean_squared_error(y_pred_xgb, y_test)]
4 errors_mae = [mean_absolute_error(y_pred, y_test), mean_absolute_error(y_rrf_pred,
5   y_test), mean_absolute_error(y_pred_xgb, y_test)]
6 r2_train = [model.score(X_train, y_train), rrf.score(X_train, y_train), xgb_model.
7   score(X_train, y_train)]
8 r2_test = [model.score(X_test, y_test), rrf.score(X_test, y_test), xgb_model.score(
9   X_test, y_test)]
df_error = pd.DataFrame({'Model name':models, 'R2 Train':r2_train, 'R2 Test':
10   r2_test, 'MSE':errors_mse, 'MAE':errors_mae})
11 df_error
```

	Model name	R2 Train	R2 Test	MSE	MAE
0	Multiple Linear/Lasso Regression	0.805799	0.807700	4.430180	1.214303
1	Random Forest	0.903110	0.767744	5.350690	1.472494
2	XGB	0.867596	0.774986	5.183843	1.406578

```
1 plt.figure(figsize=(30,8))
2 plt.subplot(1,4,1)
3 sns.barplot(df_error['Model name'], df_error['R2 Train'])
4 plt.grid()
5
6 plt.subplot(1,4,2)
7 sns.barplot(df_error['Model name'], df_error['R2 Test'])
8 plt.grid()
9
10 plt.subplot(1,4,3)
11 sns.barplot(df_error['Model name'], df_error['MSE'])
12 plt.grid()
13
14 plt.subplot(1,4,4)
15 sns.barplot(df_error['Model name'], df_error['MAE'])
16
17 plt.grid()
18 plt.legend()
19 plt.show()
```





### 10.5.2 New prediction value

I'll create a new value to make the prediction, then compare 3 prediction values together:

```
1 X_pred = [[9,13,3,1]]      #modify the input values here  (G1, G2, Medu, higher  
2   respectively)  
3 y_pred_ = [model.predict(X_pred), rrf.predict(X_pred), xgb_model.predict(X_pred)]  
4  
5 models = ["Multiple Linear/Lasso Regression", "Random Forest", "XGB"]  
6 df_pred = pd.DataFrame({'Model name':models, 'Predicted Values':y_pred_})  
6 df_pred
```

	Model name	Predicted Values
0	Multiple Linear/Lasso Regression	[12.145273062565975]
1	Random Forest	[12.99]
2	XGB	[13.0274105]

## 11 Application

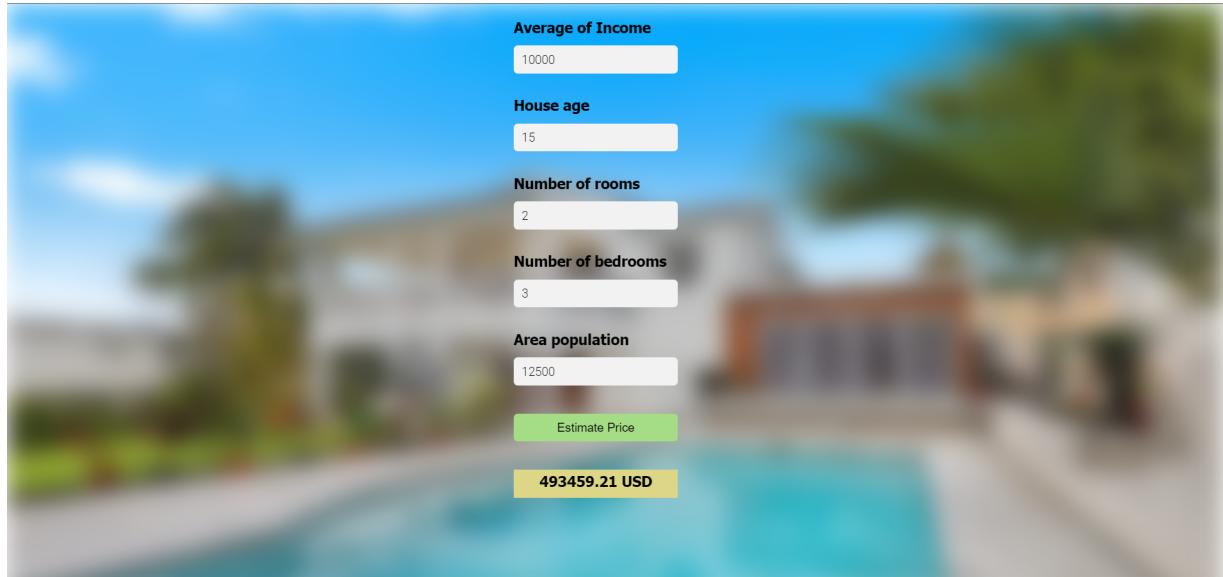
In this section, we are going to review about a web application that we have made for one simple dataset, which is about House price prediction.

Source code: [github](#)

### 11.1 Usage

Since real estate industry has always been a hot topic, we developed this application with a desire that it can support people to have an overview or a prediction about an area's house price.

Let's begin with its GUI:



There are five input fields in our application, including average of income (average income in an area), house age, number of rooms, number of bedrooms and area population. Using these inputs, our application will return for the users an estimate price, which is the predicted house price.

### 11.2 Evaluation

As this application is just a prototype for illustrating our model, its UI is now quite simple and our model does not always give the best prediction for some cases. In the future, we will try to improve it for a better version.