1) As the data are stored in Datafile, thus need load into Jupyter notebook as dataframe first before converting to numpy.
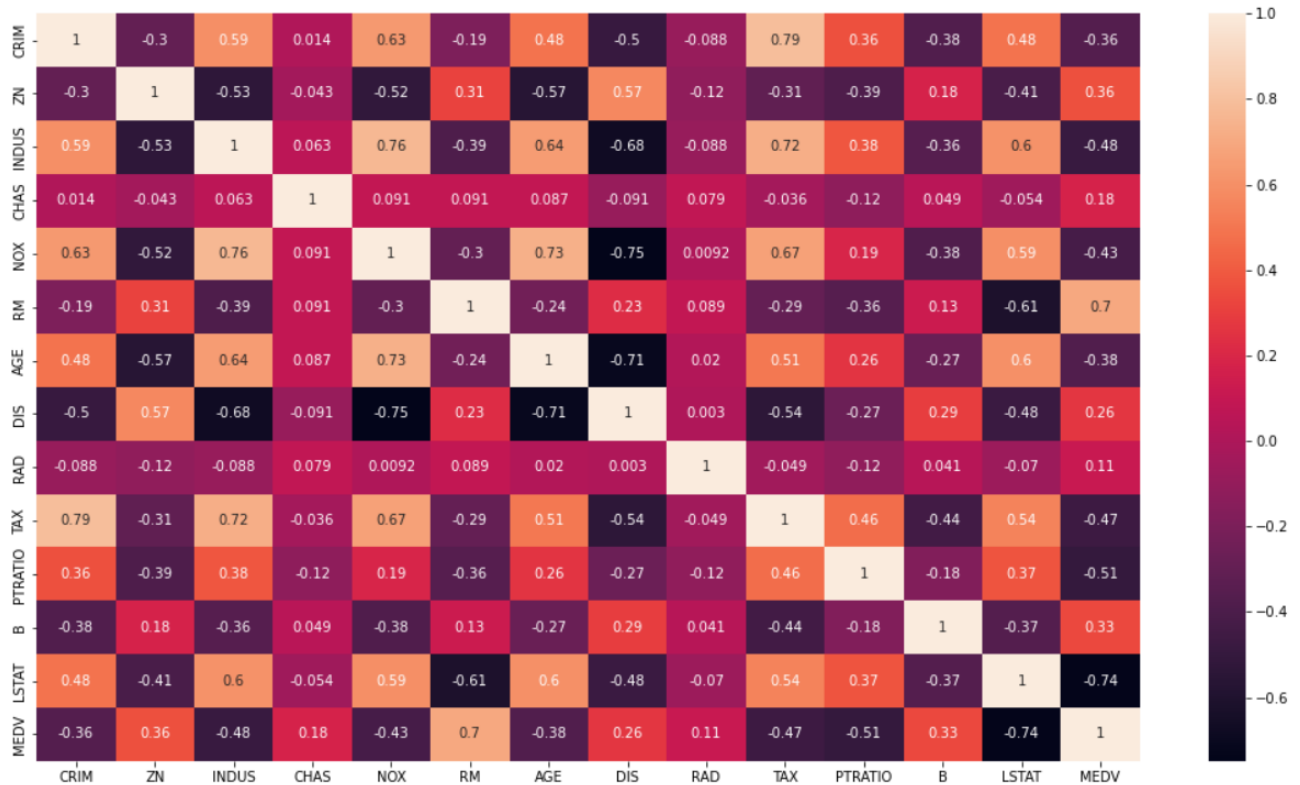
   There are no missing data for this dataset and thus, original data will be used for this case and no data modification is to be make.

2) linear regression equation:

$$y = b_0 + b_1x_1 + b_2x_2 + b_3x_3 + b_4x_4 + b_5x_5 + b_6x_6 + b_7x_7 + b_8x_8 + b_9x_9 + b_{10}x_{10} + b_{11}x_{11} + b_{12}x_{12} + b_{13}x_{13}$$

   As since there are 13 numerical input variables, thus 13 parameters are needed to determine the output(y), which is the predicted median value of owner-occupied homes in Boston.

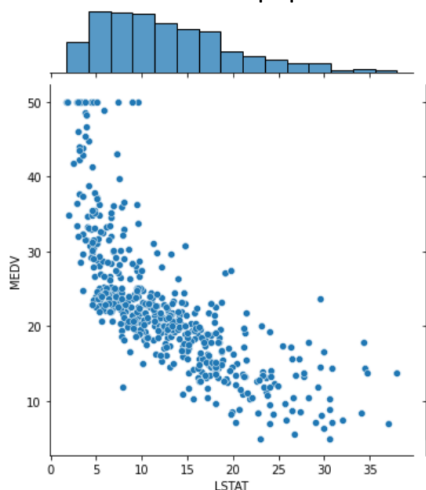3) 1$^{st}$: plot the correlation between the output vs the 13 parameters:



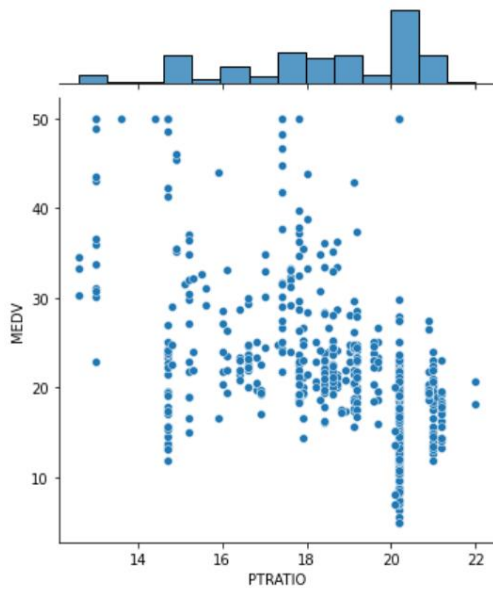2$^{nd}$: from the heatmap above, the 5 top attributes are:
- LSTAT
- RM
- PTRATIO
- TAX
- INDUS

Below are the following plots between the 5 top attributes and the output respectively:
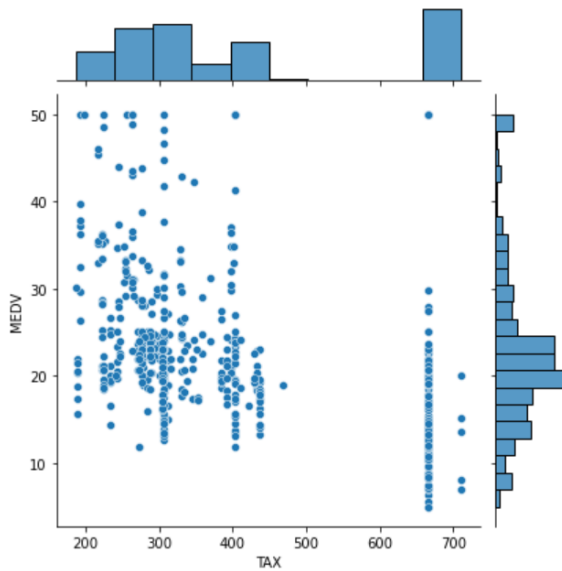
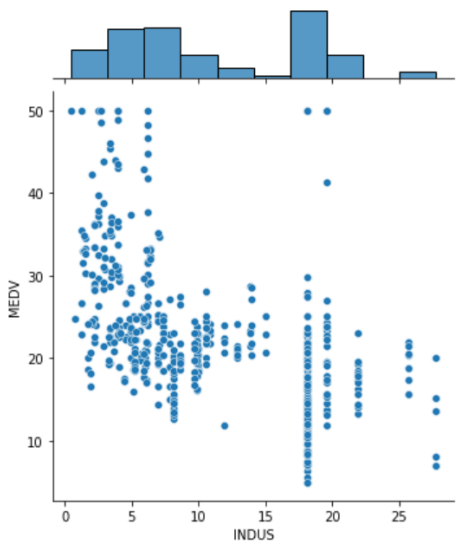⇨ % Lower status of the population vs Median value of owner-occupied homes in $1000's

⇨ pupil-teacher ratio by town vs Median value of owner-occupied homes in $1000's



⇨ full-value property-tax rate per $10,000 vs Median value of owner-occupied homes in $1000's



⇨ proportion of non-retail business acres per town vs Median value of owner-occupied homes in $1000's



4) Root Mean Square Error (RMSE) is to show how concentrated the data is around the line of best fit. The lower the RMSE, the better a given model can "fit" a dataset.
Below is the formula:

$$\text{RMSD} = \sqrt{\frac{\sum_{i=1}^{N}(x_i - \hat{x}_i)^2}{N}}$$

$\text{RMSD}$ = root-mean-square deviation
$i$ = variable i
$N$ = number of non-missing data points
$x_i$ = actual observations time series
$\hat{x}_i$ = estimated time series

R2-Score is the proportion of the variation in the dependent variable that is predictable from the independent variable(s). Generally, a higher R2-Score indicates a better fit for the model.
Below is the formula:

$$R^2 = 1 - \frac{RSS}{TSS}$$

$R^2$ = coefficient of determination
$RSS$ = sum of squares of residuals
$TSS$ = total sum of squares

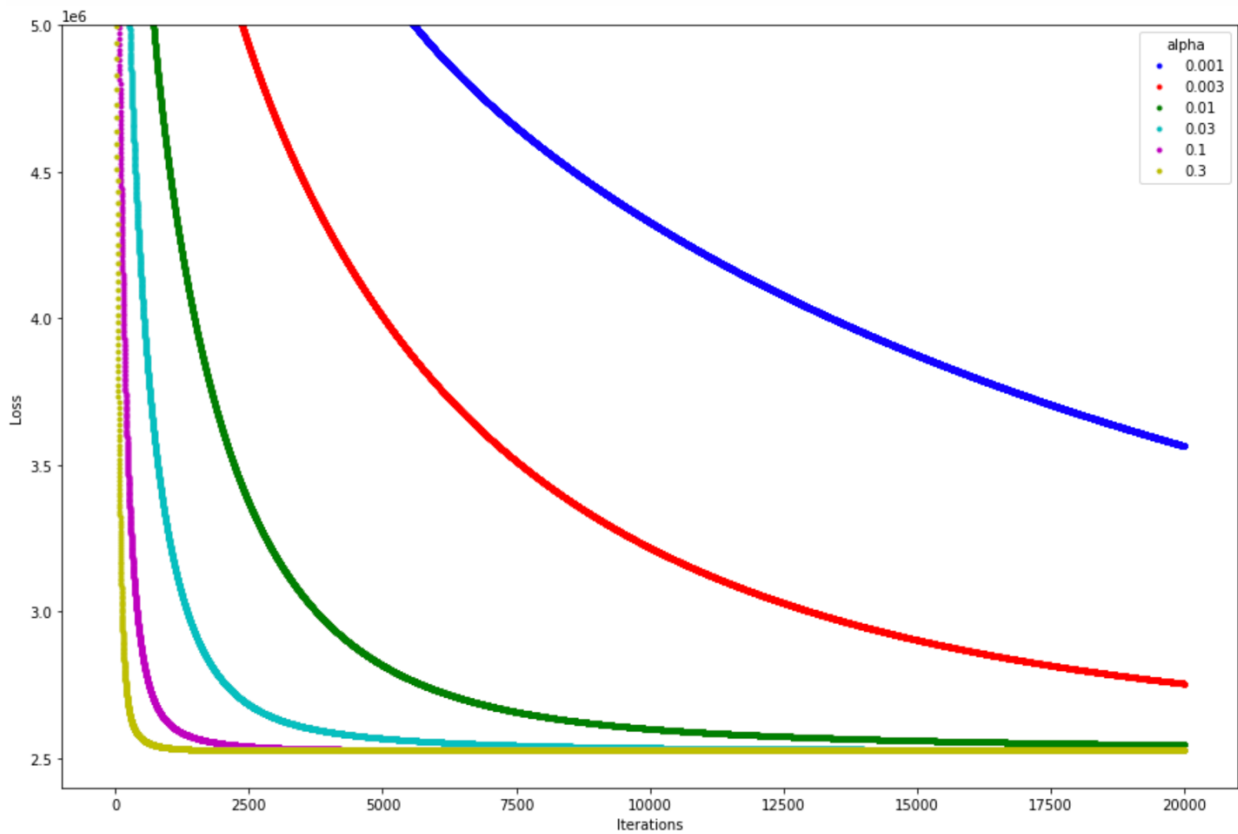Below is the python code for these 2 functions above:

```
In [ ]:
1  ### from math import sqrt
2  # Calculate root mean squared error and r2 score
3  def rmse_metric(actual, predicted):
4      sum_error = 0.0
5      for i in range(len(actual)):
6          prediction_error = predicted[i] - actual[i]
7          sum_error += (prediction_error ** 2)
8      mean_error = sum_error / float(len(actual))
9      return sqrt(mean_error)
10
11 def R2_score(actual, predicted):
12     actual_bar = actual.mean()
13     ss_tot = ((actual-actual_bar)**2).sum()
14     ss_res = ((actual-predicted)**2).sum()
15     return 1 - (ss_res/ss_tot)
```

5) Gradient Descent is to estimate the model parameters that minimize the error of the model.
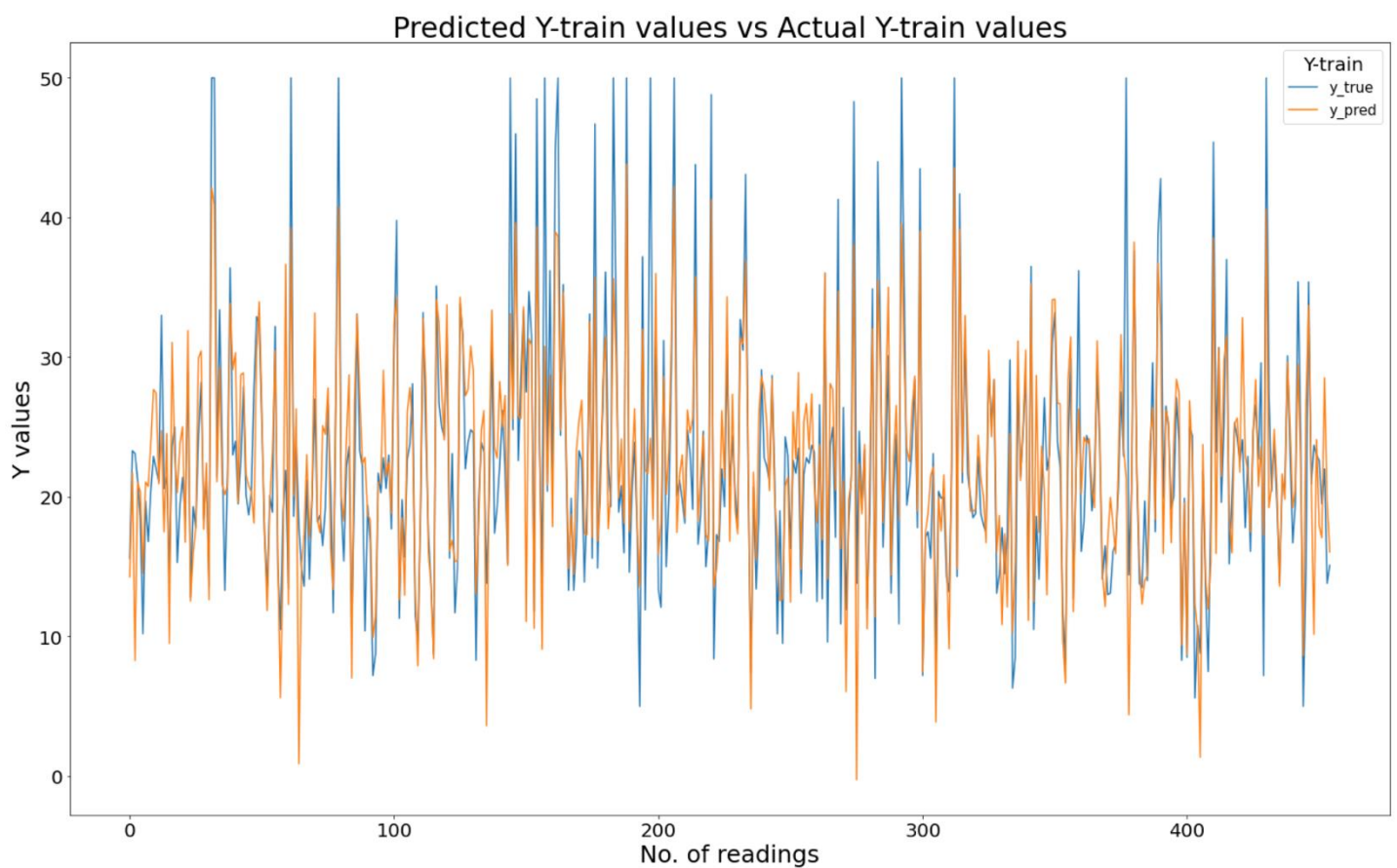Python code:

```
In [16]:
1  def gradient_descent(X,y,theta,learning_rate,iterations):
2      m = len(y)
3      n = np.shape(X)[1]
4      cost_history = np.zeros(iterations)
5      theta_history = np.zeros((iterations,n))
6      for it in range(iterations):
7          prediction = np.dot(X,theta)
8          theta = theta -(1/m)*learning_rate*(X.T.dot((prediction - y)))
9          theta_history[it,:] =theta.T
10         cost_history[it]  = cal_cost(theta,X,y)
11     return theta, cost_history, theta_history
```
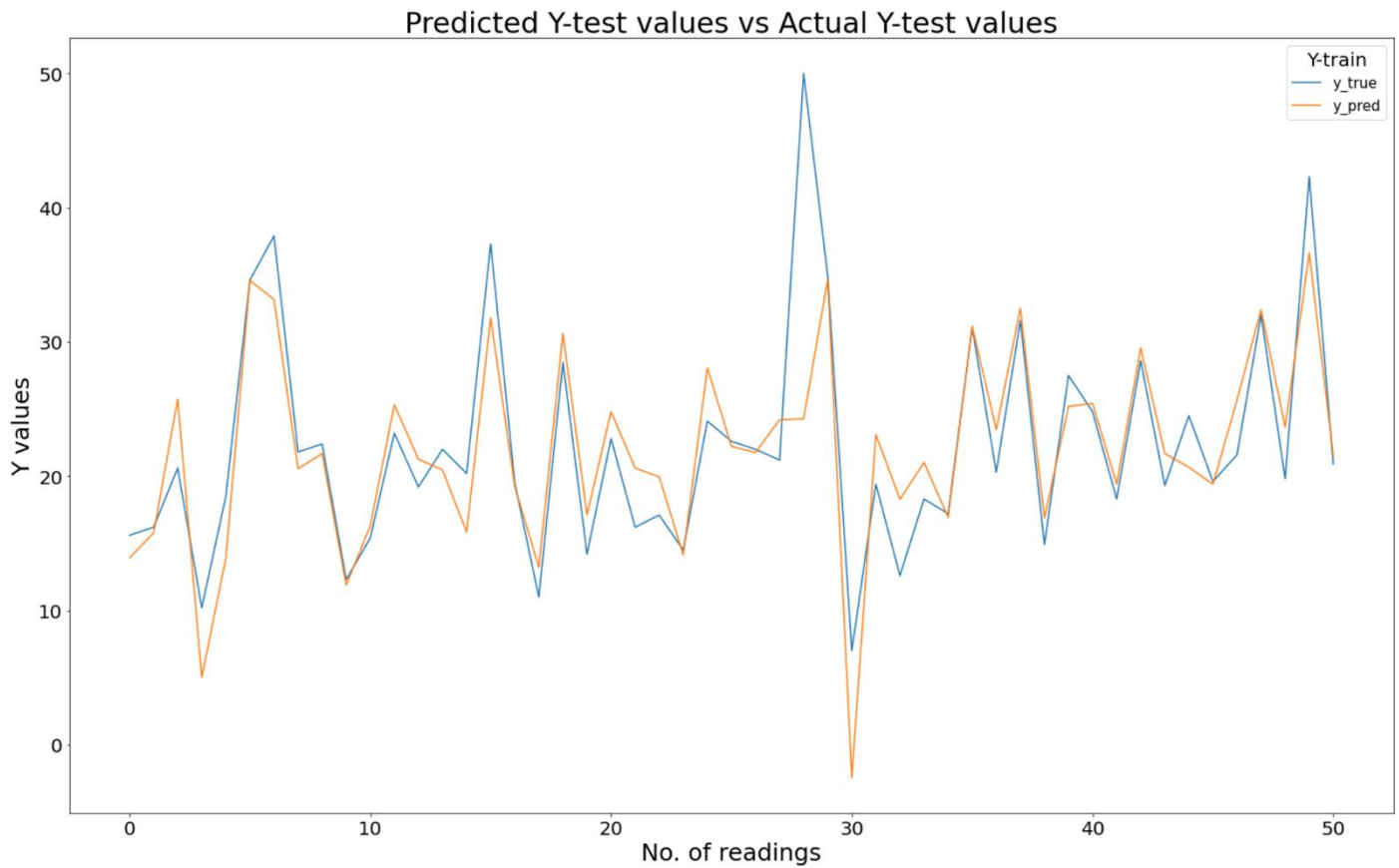
Below is the plot for different learning rate, $\alpha$ at the number of iteration=20000:



Based on the graph above, the best learning rate is $\alpha$=0.3 since it generates the least loss at the end of iterations.

Qn8)

Predicted Y-test values vs Actual Y-test values

Below is the comparison between sklearn and constructed model in terms of RMSE and R2 scores:

| | rmse | | r2_Score | |
| --- | --- | --- | --- | --- |
| | Actual | Predicted | Actual | Predicted |
| **Models** | | | | |
| **sklearn** | 4.982 | 3.73 | 0.712 | 0.795 |
| **own construct** | 4.982 | 3.73 | 0.712 | 0.795 |

Based on the graphs and rmse and r2_score, this can conclude that the constructed model is accurate as sklearn model.