# Feature Selection

Here we are performing linear regression, on a dataset with 13 features and 1 target.

| Features | Target |
|----------|--------|
|  |  |

To perfume linear regression we present 4 robust methods for feature selection:

- ## *Backward Elimination:*

This method is based on feeding all the possible features to the model at first. Then, based on the performance of the model features the worst performing features are removed. The criteria here is the pvalue. Features with pvalue > 0.05 are removed.

- pvalues for all features:



For this particular dataset 2 features have pvalue above 0.05

| Features | pvalue |
|----------|--------|
| AGE | 0.958229309205725 |
| INDUS | 0.7379887092915007 |

Based on the Backward Elimination method the selected features are: **['CRIM', 'ZN', 'CHAS', 'NOX', 'RM', 'DIS', 'RAD', 'TAX', 'PTRATIO', 'B', 'LSTAT']**

- ***Recursive Feature Elimination:***

This is a recursive method that remove features given a desired number of features. The method takes is the model 'Linear Regression' and the desired number of feature. For this case we try to optimize the number of features based on the score of the model given a number features.

For this particular case the Optimum number of features is 10
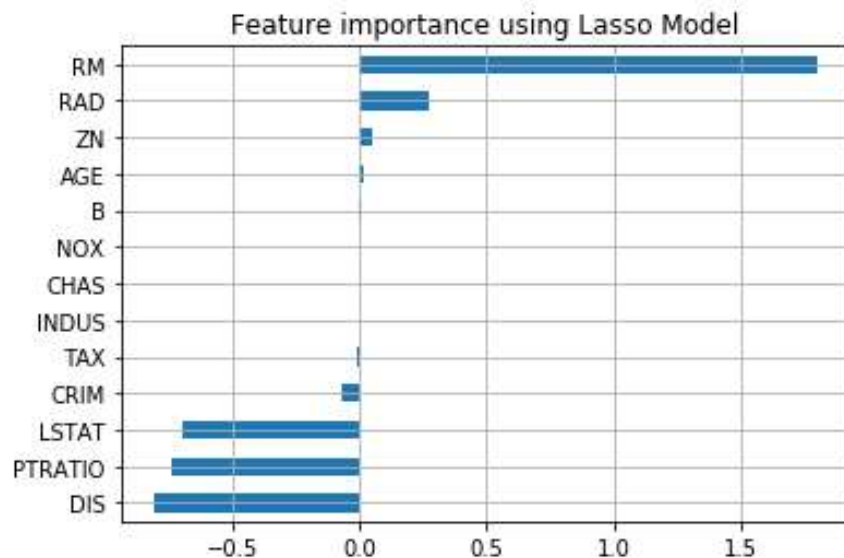
Score with 10 features is 0.663581

Based on the Recursive Feature Elimination method the selected features are: **['CRIM', 'ZN', 'INDUS', 'CHAS', 'NOX', 'RM', 'DIS', 'RAD', 'PTRATIO', 'LSTAT']**

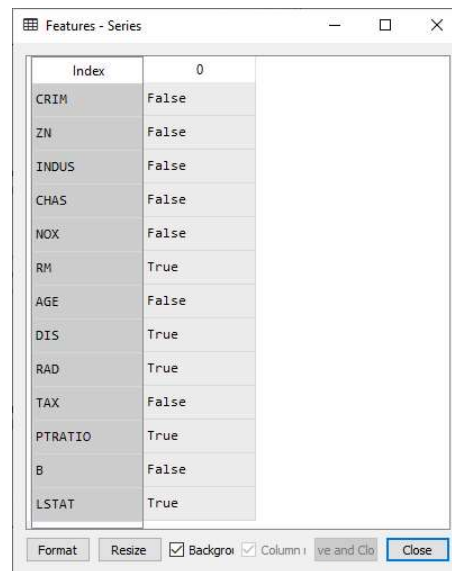- ***Embedded Method (Lasso) without threshold:***

This a regularization method that aims at extracting the features which contributes the most to the training. Here we use Lasso Regularization. Features with Lasso coefficient = 0 are removed form the model.

Based on the Embedded method (Lasso) (without threshold) the selected features are: **['CRIM', 'ZN', 'RM', 'AGE', 'DIS', 'RAD', 'TAX', 'PTRATIO', 'B', 'LSTAT']**

Lasso model picked 10 features and removed 3 features



Feature importance using Lasso Model

- *Embedded Method (Lasso) with threshold:*



| Index | 0 |
| --- | --- |
| CRIM | False |
| ZN | False |
| INDUS | False |
| CHAS | False |
| NOX | False |
| RM | True |
| AGE | False |
| DIS | True |
| RAD | True |
| TAX | False |
| PTRATIO | True |
| B | False |
| LSTAT | True |

The  model selected 5 important features and removed 8 less important features

Based on the Embedded method (Lasso) (with threshold of 0.2) the selected features are:  **['RM', 'DIS', 'RAD', 'PTRATIO', 'LSTAT']**