# RegularizationTechniques

April 5, 2023

## 1 Regularization Techniques

The purpose of Regularization techniques is to reduce the degree of overfitting that can occur in Regression models. Overfitting leads to poor ability of the model to make predictions on the new, unseen data. As we saw in the previous Regression Lessons, with a creation of extra features, such as through polynomial regression, a model can become easily overfit. To reduce the overfitting, we can regularize the model, or in other words, we can decrease its degrees of freedom. A simple way to regularize polynomial model is to reduce the number of polynomial degrees. For a linear regression model, regularization is typically achieved by constraining the weights of the model. Regularizer imposes a penalty on the size of the coefficients of the model.

Now, we will cover three types of regularizers:

- Ridge regression
- Lasso regression
- Elastic Net

Each one has its own advantages and disadvantages. Lasso will eliminate many features and reduce overfitting in your linear model. Ridge will reduce the impact of the features that are not important in predicting your target. Elastic Net combines feature elimination from Lasso and feature coefficient reduction from the Ridge model to improve your model's predictions.

The common features of all these regularizers include using cross-validation to select hyperparameters and applying data normalization to improve the performance.

**Objectives**

- Understand the advantages and disadvantages of Ridge, Lasso and Elastic Net Regressions
- Apply Ridge, Lasso and Elastic Net Regressions
- Perform hyperparameters Grid Search on a model using validation data

```python
[1]: # Surpress warnings:
     def warn(*args, **kwargs):
         pass
     import warnings
     warnings.warn = warn

     import pandas as pd
     import numpy as np
```

```
import seaborn as sns
import matplotlib.pylab as plt
%matplotlib inline

from sklearn.model_selection import train_test_split, GridSearchCV
from sklearn.linear_model import LinearRegression,Ridge,Lasso,ElasticNet
from sklearn.metrics import r2_score, mean_squared_error
from sklearn.preprocessing import PolynomialFeatures, scale, StandardScaler,␣
  ↪MinMaxScaler

from sklearn.feature_selection import SelectKBest, f_regression
from sklearn.pipeline import Pipeline
from sklearn.decomposition import PCA
```

First, let's define some functions that will help us in the future analysis. Below function will calculate the $R^2$ on each feature given the input of the model.

```
[2]: def get_R2_features(model,test=True):
         #X: global
         features=list(X)
         features.remove("three")

         R_2_train=[]
         R_2_test=[]

         for feature in features:
             model.fit(X_train[[feature]],y_train)

             R_2_test.append(model.score(X_test[[feature]],y_test))
             R_2_train.append(model.score(X_train[[feature]],y_train))

         plt.bar(features,R_2_train,label="Train")
         plt.bar(features,R_2_test,label="Test")
         plt.xticks(rotation=90)
         plt.ylabel("$R^2$")
         plt.legend()
         plt.show()
         print("Training R^2 mean value {} Testing R^2 mean value {} ".format(str(np.
  ↪mean(R_2_train)),str(np.mean(R_2_test))) )
         print("Training R^2 max value {} Testing R^2 max value {} ".format(str(np.
  ↪max(R_2_train)),str(np.max(R_2_test))) )
```

Below function will plot the estimated coefficients for each feature and find $R^2$ on training and testing sets.

```
[3]: def plot_coef(X,model,name=None):
         plt.bar(X.columns[2:],abs(model.coef_[2:]))
         plt.xticks(rotation=90)
         plt.ylabel("$coefficients$")
         plt.title(name)
         plt.show()
         print("R^2 on training  data ",model.score(X_train, y_train))
         print("R^2 on testing data ",model.score(X_test,y_test))
```

Below function plots the distribution of two inputs.

```
[4]: def  plot_dis(y,yhat):

         plt.figure()
         ax1 = sns.distplot(y, hist=False, color="r", label="Actual Value")
         sns.distplot(yhat, hist=False, color="b", label="Fitted Values" , ax=ax1)
         plt.legend()

         plt.title('Actual vs Fitted Values')
         plt.xlabel('Price (in dollars)')
         plt.ylabel('Proportion of Cars')

         plt.show()
         plt.close()
```

```
[5]: # let's load the data
     data = pd.read_csv('https://cf-courses-data.s3.us.cloud-object-storage.
       ↪appdomain.cloud/IBM-ML240EN-SkillsNetwork/labs/encoded_car_data.csv')
     data.head()
```

```
[5]:    diesel  gas  std  turbo  convertible  hardtop  hatchback  sedan  wagon  \
     0     0.0  1.0  1.0    0.0          1.0      0.0        0.0    0.0    0.0
     1     0.0  1.0  1.0    0.0          1.0      0.0        0.0    0.0    0.0
     2     0.0  1.0  1.0    0.0          0.0      0.0        1.0    0.0    0.0
     3     0.0  1.0  1.0    0.0          0.0      0.0        0.0    1.0    0.0
     4     0.0  1.0  1.0    0.0          0.0      0.0        0.0    1.0    0.0

        4wd  …  wheelbase  curbweight  enginesize  boreratio  horsepower  \
     0  0.0  …       88.6      2548.0       130.0       3.47       111.0
     1  0.0  …       88.6      2548.0       130.0       3.47       111.0
     2  0.0  …       94.5      2823.0       152.0       2.68       154.0
     3  0.0  …       99.8      2337.0       109.0       3.19       102.0
     4  1.0  …       99.4      2824.0       136.0       3.19       115.0

        carlength  carwidth  citympg  highwaympg    price
     0      168.8      64.1     21.0        27.0  13495.0
```

```
1        168.8        64.1        21.0        27.0  16500.0
2        171.2        65.5        19.0        26.0  16500.0
3        176.6        66.2        24.0        30.0  13950.0
4        176.6        66.4        18.0        22.0  17450.0

[5 rows x 36 columns]
```

[6]: `data.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 205 entries, 0 to 204
Data columns (total 36 columns):
 #   Column       Non-Null Count  Dtype
---  ------       --------------  -----
 0   diesel       205 non-null    float64
 1   gas          205 non-null    float64
 2   std          205 non-null    float64
 3   turbo        205 non-null    float64
 4   convertible  205 non-null    float64
 5   hardtop      205 non-null    float64
 6   hatchback    205 non-null    float64
 7   sedan        205 non-null    float64
 8   wagon        205 non-null    float64
 9   4wd          205 non-null    float64
 10  fwd          205 non-null    float64
 11  rwd          205 non-null    float64
 12  dohc         205 non-null    float64
 13  dohcv        205 non-null    float64
 14  l            205 non-null    float64
 15  ohc          205 non-null    float64
 16  ohcf         205 non-null    float64
 17  ohcv         205 non-null    float64
 18  rotor        205 non-null    float64
 19  eight        205 non-null    float64
 20  five         205 non-null    float64
 21  four         205 non-null    float64
 22  six          205 non-null    float64
 23  three        205 non-null    float64
 24  twelve       205 non-null    float64
 25  two          205 non-null    float64
 26  wheelbase    205 non-null    float64
 27  curbweight   205 non-null    float64
 28  enginesize   205 non-null    float64
 29  boreratio    205 non-null    float64
 30  horsepower   205 non-null    float64
 31  carlength    205 non-null    float64
 32  carwidth     205 non-null    float64
 33  citympg      205 non-null    float64
```

```
   34   highwaympg      205 non-null      float64
   35   price           205 non-null      float64
dtypes: float64(36)
memory usage: 57.8 KB
```

```python
[8]:  # data preparation: let's split the data into X features and y target
      X = data.drop('price', axis=1)
      y = data.price

      X_train, X_test, y_train, y_test = train_test_split(X,y, test_size=0.1,␣
        ↪random_state=42)
      print("number of test samples :", X_test.shape[0])
      print("number of training samples:",X_train.shape[0])
```

```
number of test samples : 21
number of training samples: 184
```

Linear Regression

In linear regression we are trying to find the value of $\mathbf{w}$ that minimizes the Mean Squared Error (MSE), we can represent this using the following expression:

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}}||\mathbf{y} - \mathbf{X}\mathbf{w}||^2\_2$$

Where $\mathbf{y}$ is the target, $\mathbf{X}$ is the training set and $\mathbf{w}$ is the parameter weights. The resulting $\hat{\mathbf{w}}$ is the best value to minimize the MSE, i.e., the distance between the target $\mathbf{y}$ and the estimate $\mathbf{X}\mathbf{w}$. We do this by fitting the model.

```python
[9]:  # Let's create a LinearRegression object, called lm
      lm = LinearRegression()
      # let's fit the model with multiple features on our X_train and y_train data.
      lm.fit(X_train, y_train)
```

```
[9]:  LinearRegression()
```

```python
[10]:  # apply predict() function on the testing data set.
       predicted = lm.predict(X_test)
```

```python
[11]:  # Let's calculate the R^2 on both, training and testing data sets
       print("R^2 on training  data ",lm.score(X_train, y_train))
       print("R^2 on testing data ",lm.score(X_test,y_test))
```
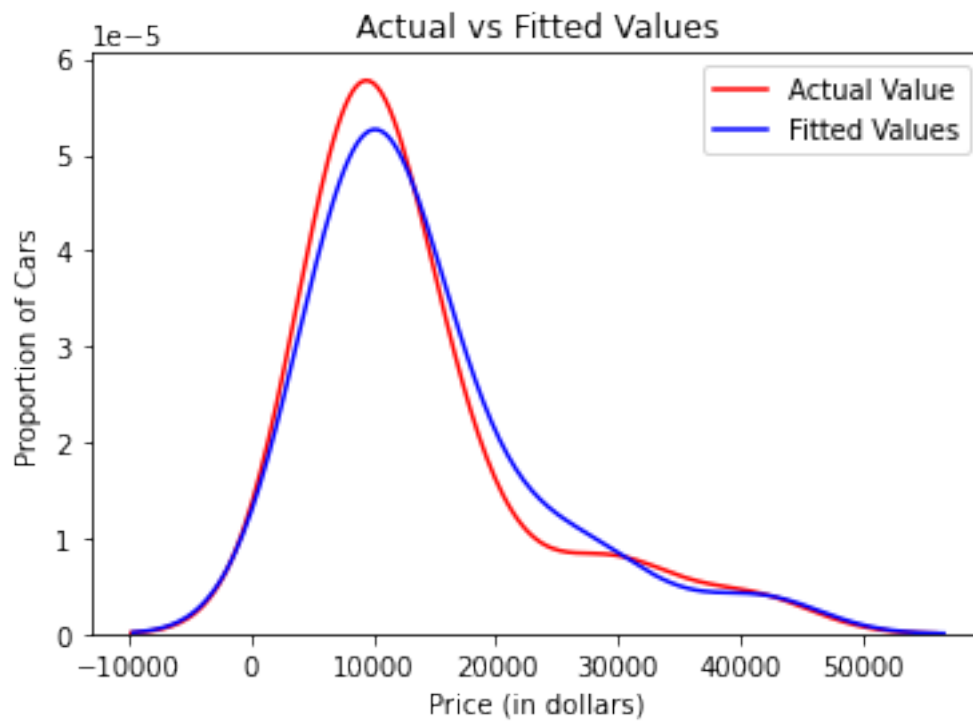
```
R^2 on training  data  0.9092101381197337
R^2 on testing data  0.9472499250320893
```
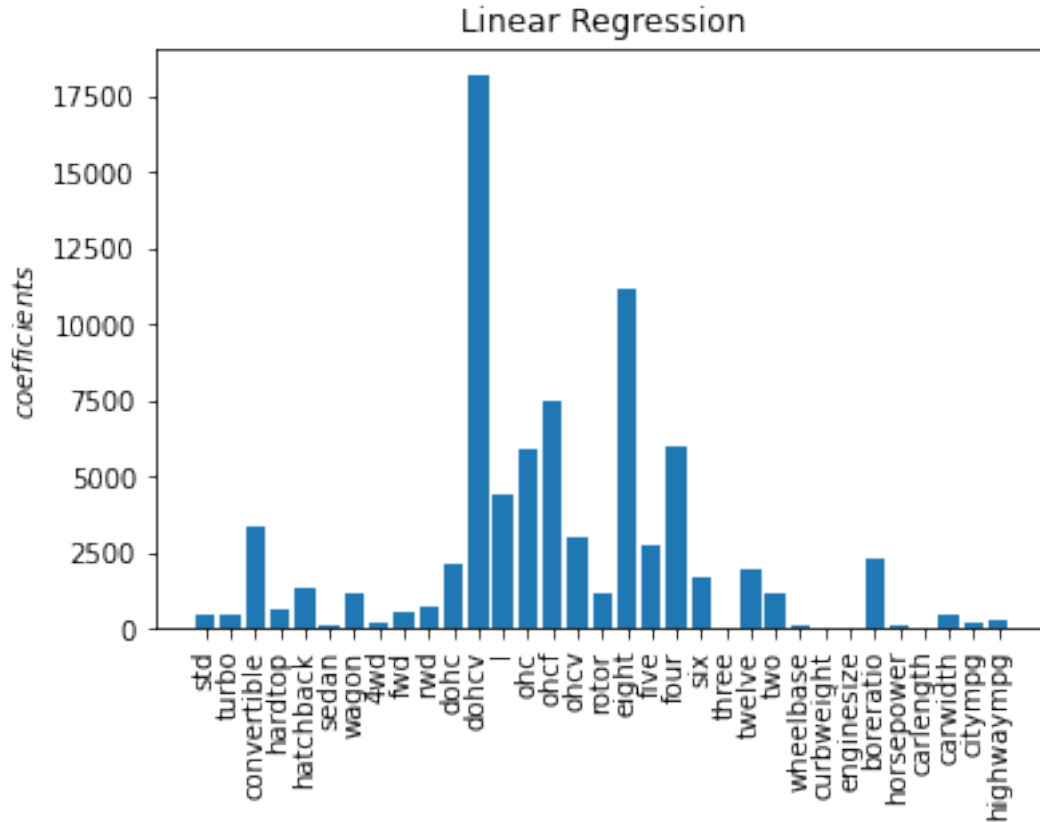
```python
[12]:  plot_dis(y_test,predicted)
```

Actual vs Fitted Values

We can view the estimated coefficients for the linear regression problem and drop the top two coefficients, as they are two large.

[13]: ```
plot_coef(X,lm,name="Linear Regression")
```

```
R^2 on training  data  0.9092101381197337
R^2 on testing data   0.9472499250320893
```

---

Ridge Regression

Let's review the Ridge Regression. Ridge Regression makes the prior assumption that our coefficients are normally distributed around zero. A regularization term, alpha, is added to the cost function. This forces the learning algorithm to not only fit the data but also keep the model weights as small as possible. The variance of the distribution is inversely proportional to the parameter alpha. This is also called the L2 regularizer , as it adds a L2 penalty to the minimization term, as shown here:

$$\hat{\mathbf{w}} = \operatorname*{argmin}_{\mathbf{w}} ||\mathbf{y} - \mathbf{X}\mathbf{w}||^2\_2 + \alpha ||\mathbf{w}||\_2$$

We minimize the MSE, but we also penalize large weights by including their magnitude $||\mathbf{w}||\_2$ in the minimization term. This additional minimization term makes the model less susceptible to noise and makes the weights smaller. Alpha controls the takeoff between MSE and penalization or regularization term and is chosen via cross-validation.

Let's see how the parameter alpha changes the model. Note, here our test data will be used as

validation data. Also, the regularization term should only be added to the cost function during the training.

Let's create a Ridge Regression object, setting the regularization parameter (alpha) to 0.01.

```
[14]: rr = Ridge(alpha=0.01)
      rr
```

```
[14]: Ridge(alpha=0.01)
```

```
[15]: # Like regular regression, you can fit the model using the fit() method
      rr.fit(X_train, y_train)
```

```
[15]: Ridge(alpha=0.01)
```

```
[16]: # obtain a prediction
      rr.predict(X_test)
```

```
[16]: array([30178.77172992, 22179.93145434, 11229.58960483, 11790.40337149,
             26348.13785546,  5439.13547145,  9054.02541015,  7265.36558563,
             10591.48456189, 10390.82134687, 17471.12024994,  7010.47430228,
             16547.06078383, 10468.27937016, 41540.38102791,  5390.34659152,
              5109.76377302, 15373.1930261 , 10703.56615831, 11448.33242702,
             10565.49269055])
```

```
[17]: # calculate the  2 on the training and testing data.
      print("R^2 on training  data ",rr.score(X_train, y_train))
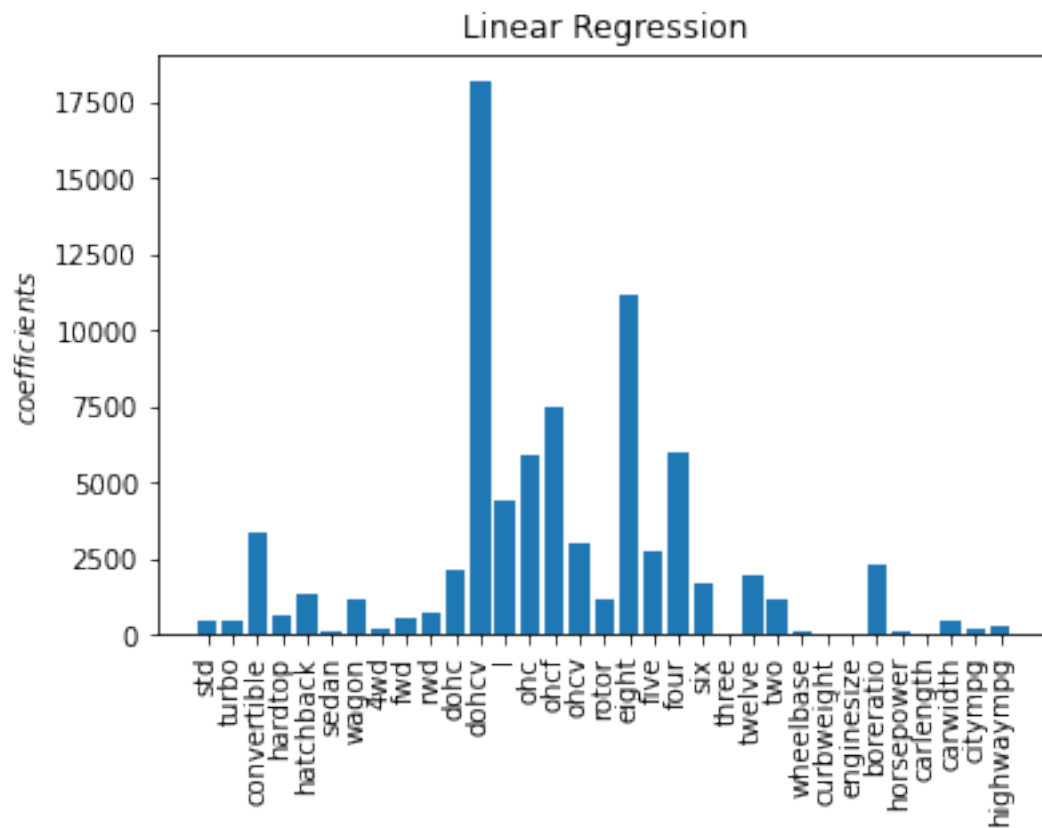      print("R^2 on testing data ",rr.score(X_test,y_test))
```

```
R^2 on training  data  0.9091956531801182
R^2 on testing data  0.9478784615596494
```

Now let's compare the Ridge Regression and the Linear Regression models. The results on the $R^2$ are about the same, and the coefficients seem to be smaller.

```
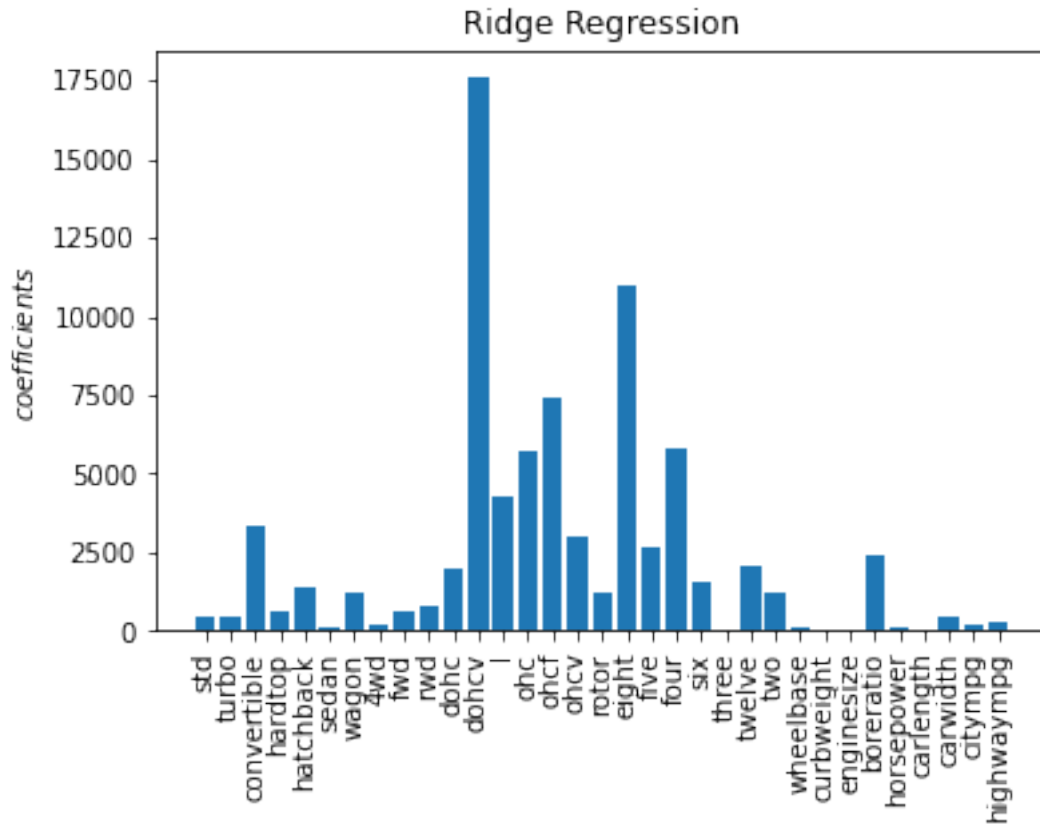[18]: plot_coef(X,lm,name="Linear Regression")
      plot_coef(X,rr,name="Ridge Regression")
```

Linear Regression

R^2 on training  data   0.9092101381197337
R^2 on testing data   0.9472499250320893

## Ridge Regression



```
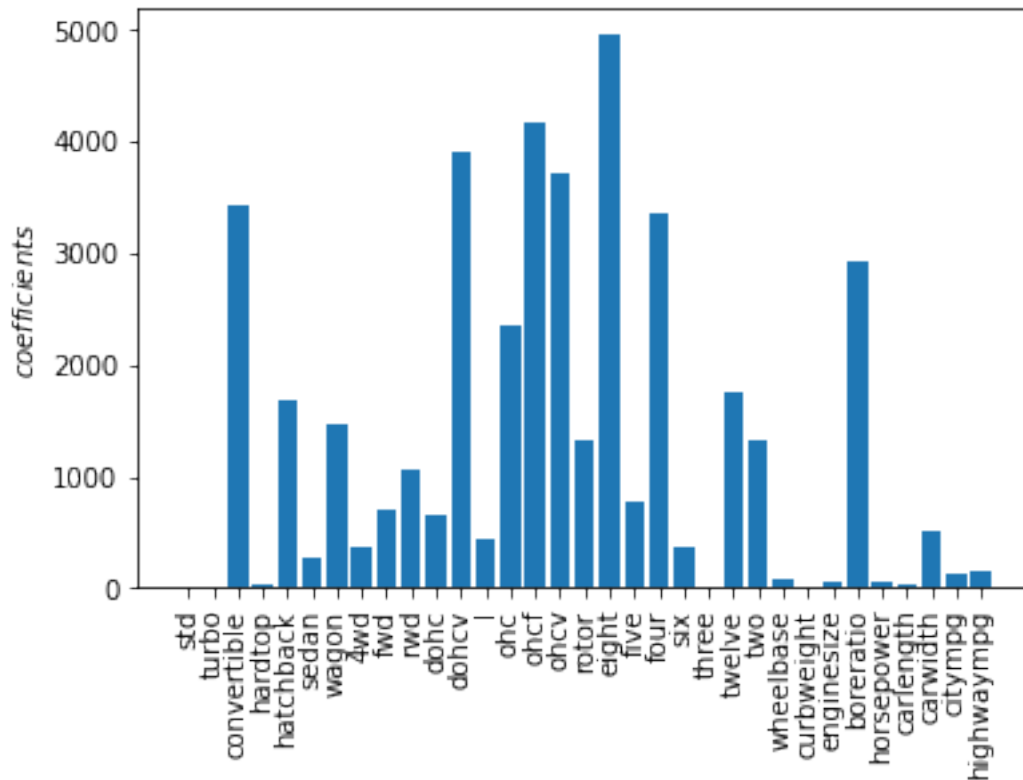R^2 on training  data  0.9091956531801182
R^2 on testing data   0.9478784615596494
```

If we increase alpha, the coefficients get smaller, but the results are not as good as our previous value of alpha.

```
[19]:  rr = Ridge(alpha=1)
       rr.fit(X_train, y_train)
       plot_coef(X,rr)
```

```
R^2 on training  data  0.8991374778636106
R^2 on testing data  0.9446031107273962
```

In general, we see that if we increase alpha, the coefficients get smaller, but the model performance relationship gets more complex. As a result, we use the validation data to select a value for alpha. Here, we plot the coefficients and R^2 of the test data on the vertical axes and alpha on the horizontal axis, as well the R^2 using the test data.

```
[20]: alphas = [0.00001,0.0001,0.001,0.01,0.1,1,10,100]
      R_2=[]
      coefs = []
      for alpha in alphas:
          ridge = Ridge(alpha=alpha)
          ridge.fit(X_train, y_train)
          coefs.append(abs(ridge.coef_))
          R_2.append(ridge.score(X_test,y_test))


      ax = plt.gca()
      ax.plot(alphas, coefs)
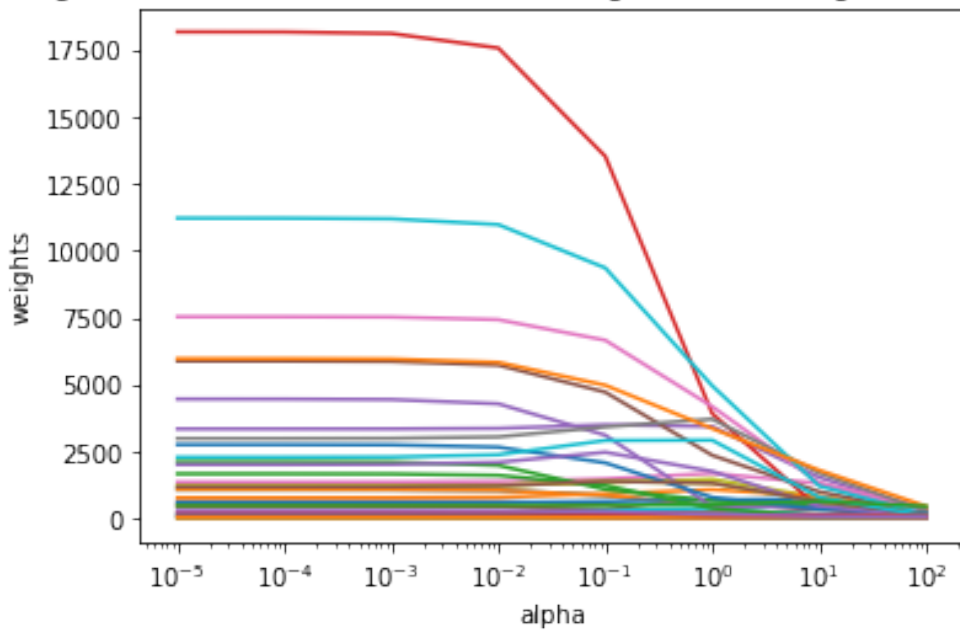      ax.set_xscale("log")
      plt.xlabel("alpha")
```

11

```
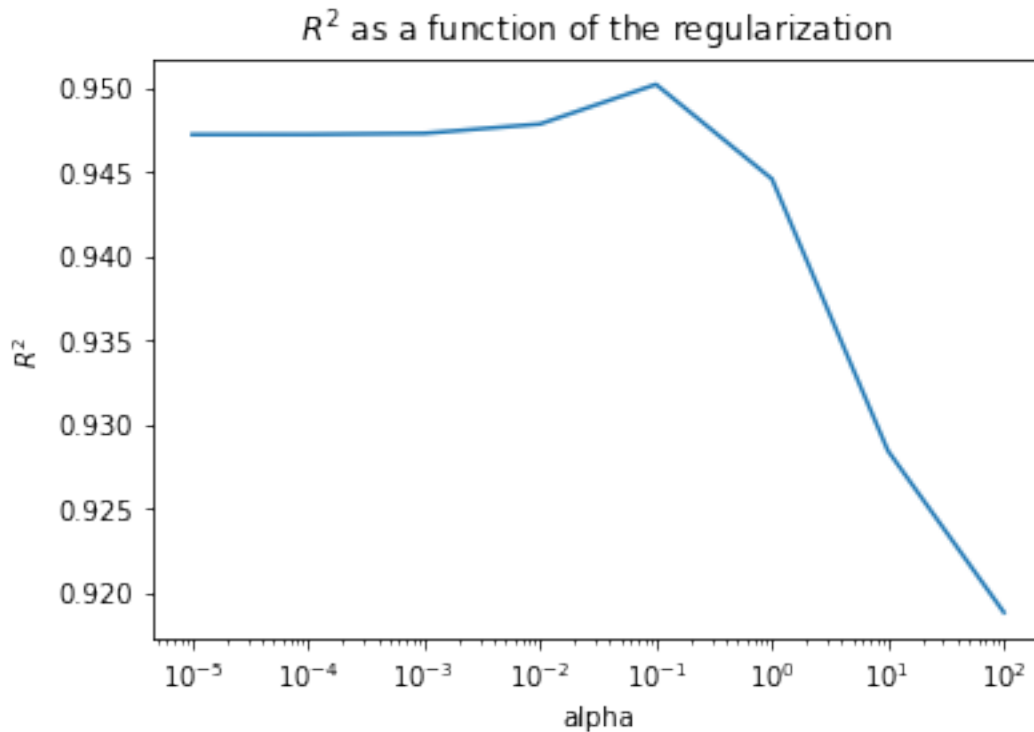plt.ylabel("weights")
plt.title("Ridge coefficients as a function of the regularization␣
  ↪(regularization path)")
plt.show()


ax = plt.gca()
ax.plot(alphas, R_2)
ax.set_xscale("log")
plt.xlabel("alpha")
plt.ylabel("$R^2$")
plt.title("$R^2$ as a function of the regularization")
plt.show()
```



Ridge coefficients as a function of the regularization (regularization path)

$R^2$ as a function of the regularization

As we increase alpha, the coefficients get smaller but the R^2 peaks when alpha is 1.

Exercise 1

Plot the MSE as a function of alpha. What pattern do you notice?

```
[60]: alphas = [0.00001,0.0001,0.001,0.01,0.1,1,10]
MEAN_SQE=[]

for alpha in alphas:
    ridge = Ridge(alpha=alpha)
    ridge.fit(X_train, y_train)
    MEAN_SQE.append(mean_squared_error(ridge.predict(X_test),y_test))

ax = plt.gca()
ax.plot(alphas, MEAN_SQE)
ax.set_xscale("log")
plt.xlabel("alpha")
plt.ylabel("MSE")
plt.title("MSE as a function of the regularization")
plt.show()
```

13

$\rightarrow$ A small alpha leads to over-fitting but as alpha gets larger the MSE decreases. When alpha gets too large the MSE increases leading to underfitting. The optimal point seems to be in the middle.

---

Pipeline

We can also create a Pipeline object and apply a set of transforms sequentially. Then, we can apply Polynomial Features, perform data standardization then apply Ridge regression. Data Pipelines simplify the steps of processing the data. We use the module `Pipeline` to create a pipeline. We also use `StandardScaler` step in our pipeline. Scaling our data is necessary step in Ridge regression as it will penalize features with a large magnitude.

```python
[22]: # create a pipeline object
Input=[ ('polynomial',
  PolynomialFeatures(include_bias=False,degree=2)),('ss',StandardScaler() ),
  ('model',Ridge(alpha=1))]
pipe = Pipeline(Input)

# fit
pipe.fit(X_train, y_train)
```

```
[22]: Pipeline(steps=[('polynomial', PolynomialFeatures(include_bias=False)),
                ('ss', StandardScaler()), ('model', Ridge(alpha=1))])
```

Looking for hyperparameters can get difficult with loops. The problem will get worse as we add

14

more transforms such as polynomial transform. Therefore, we can use `GridSearchCV` to make things simpler.

---

GridSearchCV

To search for the best combination of hyperparameters we can create a `GridSearchCV()` function as a dictionary of parameter values. The parameters of pipelines can be set by using the name of the key, separated by "___", then the parameter. Here, we look for different polynomial degrees and different values of alpha.

```
[23]: param_grid = {
          "polynomial__degree": [1,2,3,4],
          "model__alpha":[0.0001,0.001,0.01,0.1,1,10]
      }
```

Keys of the dictionary are the model "key name ___" followed by the parameter as an attribute.

polynomial___degree: is the degree of the polynomial; in this case 1, 2, 3, 4 and 5. model___alpha : Regularization strength; must be a positive float.

We create a `GridSearchCV` object and fit it. The method trains the model and the hyperparameters are selected via exhaustive search over the specified values.

```
[24]: search = GridSearchCV(pipe, param_grid, n_jobs=2)

      search.fit(X_train, y_train)
      search
```

```
[24]: GridSearchCV(estimator=Pipeline(steps=[('polynomial',
                                               PolynomialFeatures(include_bias=False)),
                                              ('ss', StandardScaler()),
                                              ('model', Ridge(alpha=1))]),
                   n_jobs=2,
                   param_grid={'model__alpha': [0.0001, 0.001, 0.01, 0.1, 1, 10],
                               'polynomial__degree': [1, 2, 3, 4]})
```

We can input the results into *pandas* `DataFrame()` as a dictionary with keys as column headers and values as columns and display the results.

```
[25]: pd.DataFrame(search.cv_results_).head()
```

```
[25]:    mean_fit_time  std_fit_time  mean_score_time  std_score_time  \
      0       0.022092      0.007494         0.001608        0.003216
      1       0.025488      0.005946         0.008694        0.007373
      2       0.129657      0.017348         0.018185        0.002850
      3       1.049331      0.031325         0.101343        0.004251
      4       0.013699      0.002070         0.007134        0.002549

         param_model__alpha param_polynomial__degree  \
```

```
0          0.0001                    1
1          0.0001                    2
2          0.0001                    3
3          0.0001                    4
4          0.001                     1
```

```
                                                params  split0_test_score  \
0  {'model__alpha': 0.0001, 'polynomial__degree': 1}           0.770265
1  {'model__alpha': 0.0001, 'polynomial__degree': 2}          -0.910593
2  {'model__alpha': 0.0001, 'polynomial__degree': 3}          -3.774944
3  {'model__alpha': 0.0001, 'polynomial__degree': 4}          -6.799155
4   {'model__alpha': 0.001, 'polynomial__degree': 1}           0.770258

   split1_test_score  split2_test_score  split3_test_score  split4_test_score  \
0           0.766681           0.844321           0.873196           0.612164
1           0.084489          -0.549646          -0.705327          -2.630419
2          -2.241587          -0.830641          -2.048391          -9.065623
3          -9.100388          -0.935754          -3.845688          -8.975313
4           0.766696           0.844330           0.873212           0.612239

   mean_test_score  std_test_score  rank_test_score
0         0.773325        0.090599                9
1        -0.942299        0.907235               21
2        -3.592237        2.892247               23
3        -5.931260        3.142215               24
4         0.773347        0.090577                8
```

There are some other useful attributes:

best_score_: mean cross-validated score of the best_estimator. best_params_dict: parameter setting that gives the best results on the hold-out data.

```python
[26]: print("best_score_: ",search.best_score_)
      print("best_params_: ",search.best_params_)
```

```
best_score_:  0.8681282282391246
best_params_:  {'model__alpha': 10, 'polynomial__degree': 2}
```

```python
[27]: # call predict() on the estimator with the best found parameters
      predict = search.predict(X_test)
      predict
```

```
[27]: array([34464.68303015, 23419.34979609,  9723.38208924, 13433.42551625,
             25239.67969507,  6532.54733115,  7468.55986929,  7687.50892226,
              9114.59988701,  9541.72323357, 15453.55023026,  6815.64031583,
             16958.71711106, 10290.1730189 , 42072.30176794,  6356.53352923,
              1536.17210277, 13705.75183086,  9758.24171738,  9210.40054167,
              9999.70716056])
```

```
[28]:  #  find the best model
       best=search.best_estimator_
       best
```

```
[28]:  Pipeline(steps=[('polynomial', PolynomialFeatures(include_bias=False)),
                       ('ss', StandardScaler()), ('model', Ridge(alpha=10))])
```

```
[29]:  predict = best.predict(X_test)
       predict
```

```
[29]:  array([34464.68303015, 23419.34979609,  9723.38208924, 13433.42551625,
              25239.67969507,  6532.54733115,  7468.55986929,  7687.50892226,
               9114.59988701,  9541.72323357, 15453.55023026,  6815.64031583,
              16958.71711106, 10290.1730189 , 42072.30176794,  6356.53352923,
               1536.17210277, 13705.75183086,  9758.24171738,  9210.40054167,
               9999.70716056])
```

```
[30]:  # now, train our model on the entire data set
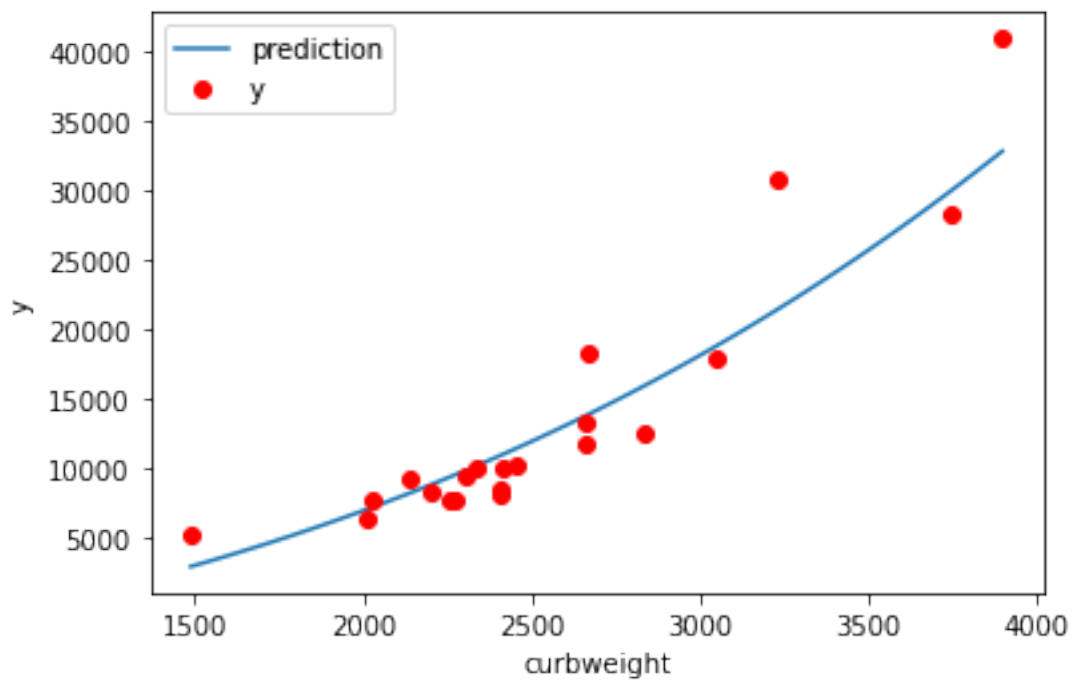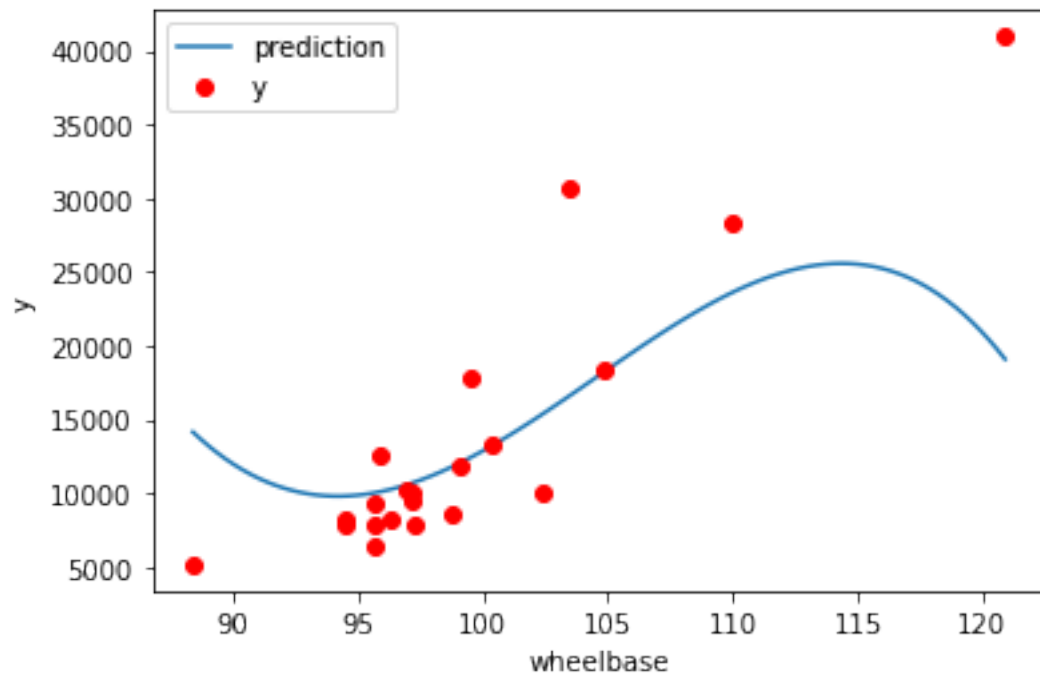       best.fit(X,y)
```

```
[30]:  Pipeline(steps=[('polynomial', PolynomialFeatures(include_bias=False)),
                       ('ss', StandardScaler()), ('model', Ridge(alpha=10))])
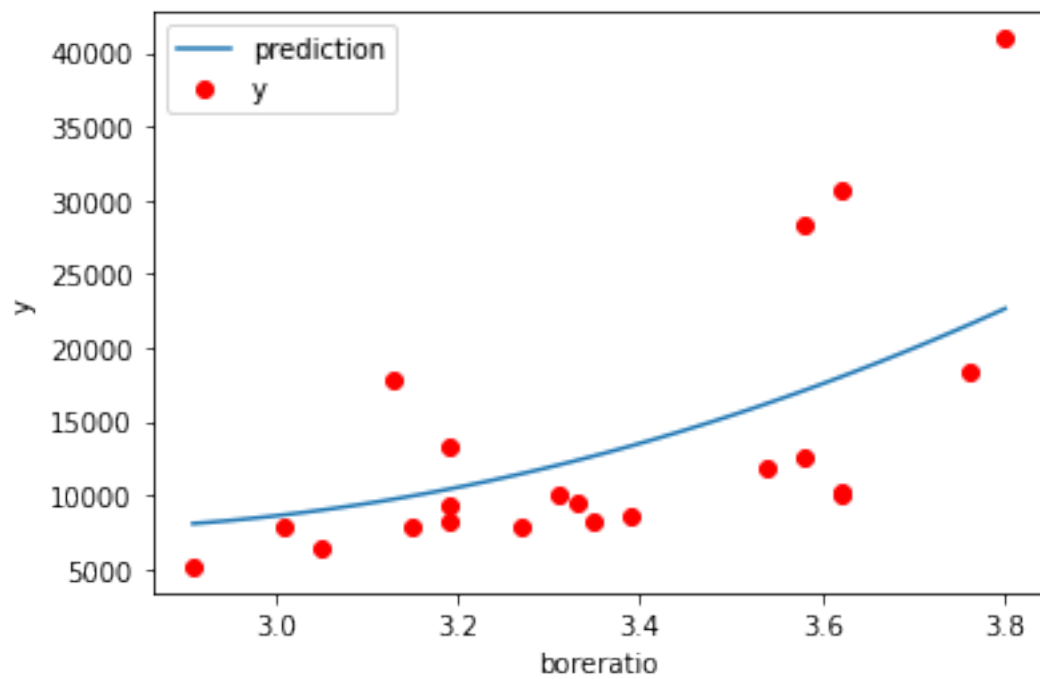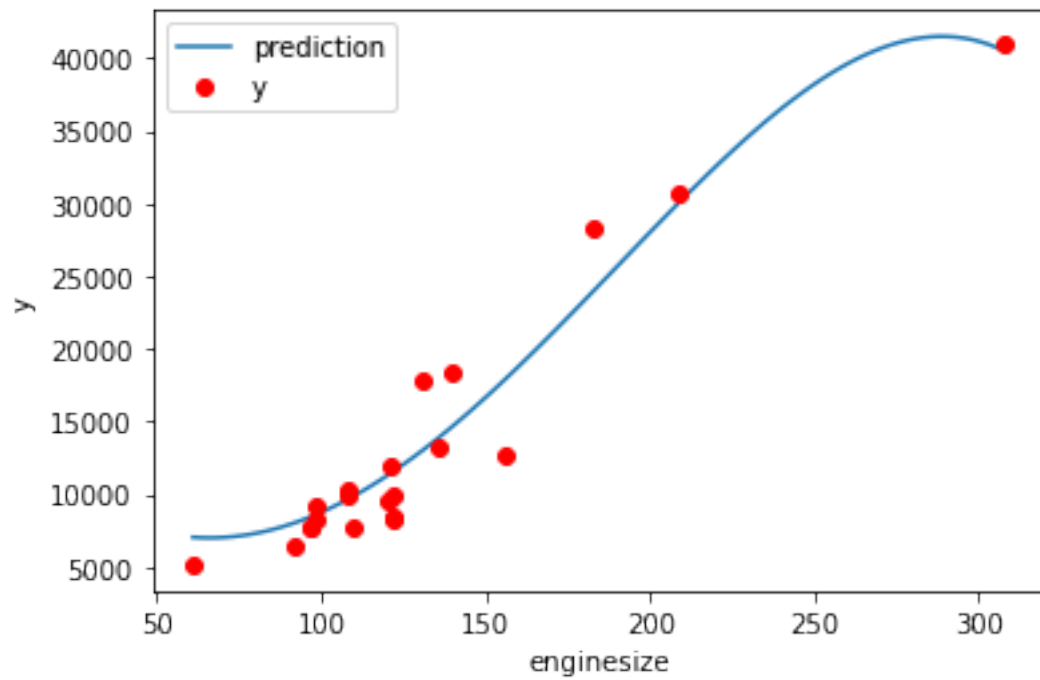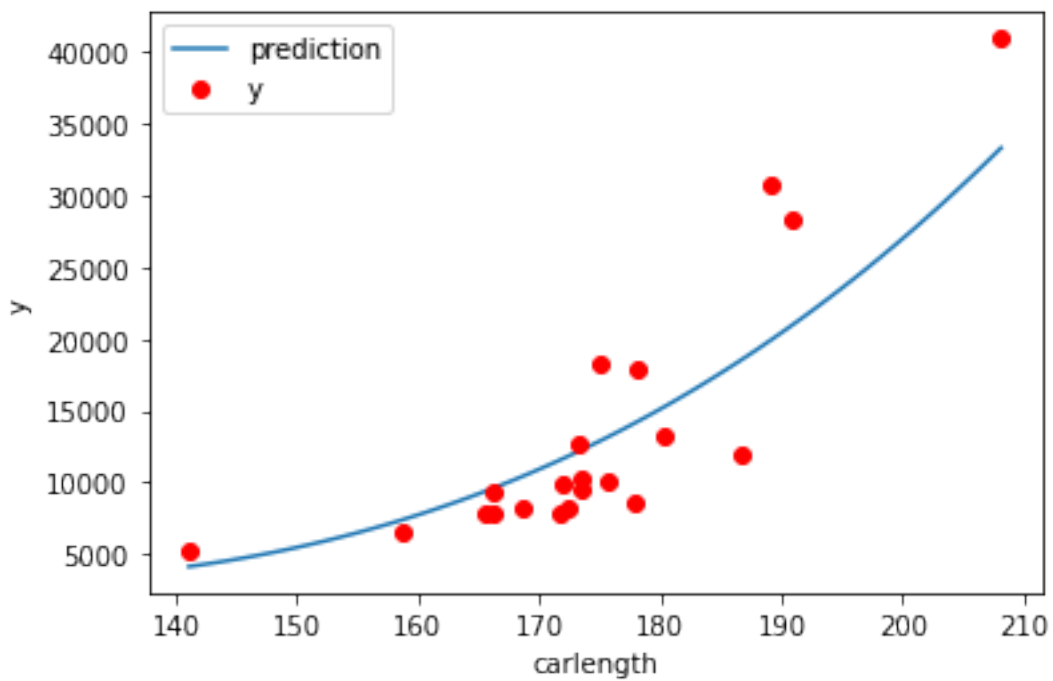```

Exercise 2

Perform grid search on the following features and plot the results by completing the following lines of code:

```
[31]:  columns=['wheelbase', 'curbweight', 'enginesize', 'boreratio', 'horsepower',
               'carlength', 'carwidth', 'citympg']

       for column in columns:
           search.fit(X_train[[column]], y_train)
           x=np.linspace(X_test[[column]].min(), X_test[[column]].max(),num=100)
           plt.plot(x,search.predict(x.reshape(-1,1)),label="prediction")
           plt.plot(X_test[column],y_test,'ro',label="y")
           plt.xlabel(column)
           plt.ylabel("y")
           plt.legend()
           plt.show()
```

Lasso Regression

In this section, let's review the Lasso (Least Absolute Shrinkage and Selection Operator) Regression. Lasso Regression makes the prior assumption that our coefficients have Laplace (double-exponential) distribution around zero. The scale parameter of the distribution is inversely proportional to the parameter alpha. The main advantage of LASSO Regression is that many coefficients are set to zero, therefore they are not required. This has many advantages, one of them is that you may not need to collect and/or store all of the features. This may save resources. For example, if the feature was some medical test, you would no longer need to perform that test. Let's see how the parameter alpha changes the model. We minimize the MSE, but we also penalize large weights by including their sum of absolute values || w || _1, symbolically:

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} ||\mathbf{y} - \mathbf{Xw}||^2 \_2 + \alpha ||\mathbf{w}||\_1$$

This regularization or penalty term makes many coefficients zero, making the model easy to understand and can also be used for feature selection. There are some drawbacks to this technique. It takes longer time to train and the solution may not be unique. Alpha controls the trade-off between MSE and penalization or regularization term and is chosen via cross-validation. Let's see how the parameter alpha changes the model. Note, as before, our test data will be used as validation data. Let's create a Ridge Regression object, setting the regularization parameter (alpha) to 0.01.

```
[32]: la = Lasso(alpha=0.1)
      la.fit(X_train,y_train)
      la
```

```
[32]: Lasso(alpha=0.1)
```

```
[33]: predicted = la.predict(X_test)
      predicted
```

```
[33]: array([30207.01664593, 22253.28978157, 11234.4293658 , 11825.51411381,
             26349.51670682,  5442.67340587,  9115.38970742,  7262.37834971,
             10602.22642602, 10426.6507622 , 17480.62216135,  6988.60683622,
             16570.42746876, 10462.5969077 , 41586.07832506,  5399.25186043,
              3317.54502062, 15367.96191917, 10698.17086097, 11451.72644684,
             10548.16668782])
```

Let's calculate the R^2 on the training and testing data and see how it performs compared to the other methods.

```
[34]: print("R^2 on training  data ",lm.score(X_train, y_train))
      print("R^2 on testing data ",lm.score(X_test,y_test))
```

```
R^2 on training  data  0.9092101381197337
R^2 on testing data  0.9472499250320893
```

If we compare the Lasso Regression to the Ridge Regression model we see that the results on the R^2 are slightly worse, but most of the coefficients are zero.

```
[35]: plot_coef(X,rr,name="Ridge Regression")
      plot_coef(X,la,name="Lasso Regression")
```



Ridge Regression

```
R^2 on training  data  0.8991374778636106
R^2 on testing data  0.9446031107273962
```

Lasso Regression

```
R^2 on training  data  0.9092098726971481
R^2 on testing data  0.9453258869077992
```

Similar to the Ridge Regression, if we increase the value of alpha, the coefficients will get smaller. Additionally, many coefficients become zero. Moreover, the model performance relationship becomes more complex. As a result, we use the validation data to select a value for alpha. Here, we plot the coefficients and R^2 of the test data on the vertical axes and alpha values on the horizontal axis.

[36]:
```python
alphas = [0.00001,0.0001,0.001,0.01,0.1,1,10,100,1000]
R_2=[]
coefs = []
for alpha in alphas:
    la=Lasso(alpha=alpha)

    la.fit(X_train, y_train)
    coefs.append(abs(la.coef_))
    R_2.append(la.score(X_test,y_test))
```

```
ax = plt.gca()
ax.plot(alphas, coefs)
ax.set_xscale("log")
plt.xlabel("alpha")
plt.ylabel("weights")
plt.title("Ridge coefficients as a function of the regularization␣
 ↪(regularization path)")
plt.show()
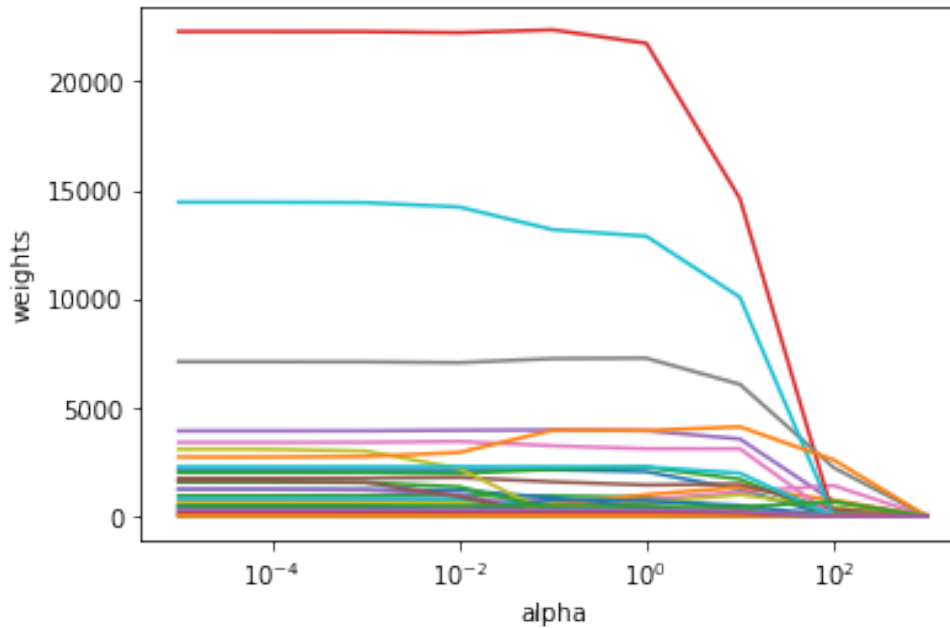

ax = plt.gca()
ax.plot(alphas, R_2)
ax.set_xscale("log")
plt.xlabel("alpha")
plt.ylabel("R^2")
plt.title("R^2 as a function of the regularization")
plt.show()
```



Ridge coefficients as a function of the regularization (regularization path)

$R^2$ as a function of the regularization

Pipeline

We can also create a Pipeline object and apply a set of transforms sequentially. Then, we can apply polynomial features, perform data standardization, then apply Lasso Regression. We also use `StandardScaler` as a step in our pipeline. Scaling your data is necessary step in LASSO Regression, as it will penalize features with a large magnitudes.

```
[37]: # start by creating a pipeline object
Input=[ ('polynomial',
  ↪PolynomialFeatures(include_bias=False,degree=2)),('ss',StandardScaler() ),
  ↪('model',Lasso(alpha=1, tol = 0.2))]
pipe = Pipeline(Input)
# fit and predict
pipe.fit(X_train, y_train)
pipe.predict(X_test)
```

```
[37]: array([31090.35502413, 23135.62504126,  9457.4650858 , 11222.64465089,
       28503.13820242,  6821.35081552,  7716.68381925,  7436.33736593,
        7805.79226801, 10151.89082851, 16262.8074276 ,  8278.75257849,
       16684.95563795, 10658.46233607, 38313.26144272,  7039.4378737 ,
        8259.35558463, 12989.37400724,  9251.96079447,  8253.92933959,
        9674.77573004])
```

```
[38]: # calculate the   2 on the training and testing data sets
      print("R^2 on training  data ",pipe.score(X_train, y_train))
      print("R^2 on testing data ",pipe.score(X_test,y_test))
```

```
R^2 on training  data  0.941181185417151
R^2 on testing data  0.953382982655632
```

As we see, some individual features perform similarly to using all the features (we removed the feature three ). Additionally, we see the smaller coefficients seem to correspond to a larger $R^2$, therefore larger coefficients correspond to overfitting.

---

GridSearchCV

To search for the best combination of hyperparameters, we can create a `GridSearchCV()` function as a dictionary of parameter values. The parameters of pipelines can be set by using the name of the key, separated by "___", then the parameter. Here, we look for different polynomial degrees and different values of alpha.

```
[39]: param_grid = {
          "polynomial__degree": [ 1, 2,3,4,5],
          "model__alpha":[0.0001,0.001,0.01,0.1,1,10]
      }
```

To search for the best combination of hyperparameters, we create a `GridSearchCV` object with a dictionary of parameter values.

```
[40]: search = GridSearchCV(pipe, param_grid, n_jobs=2)
      search.fit(X_train, y_train)
```

```
[40]: GridSearchCV(estimator=Pipeline(steps=[('polynomial',
                                               PolynomialFeatures(include_bias=False)),
                                              ('ss', StandardScaler()),
                                              ('model', Lasso(alpha=1, tol=0.2))]),
                   n_jobs=2,
                   param_grid={'model__alpha': [0.0001, 0.001, 0.01, 0.1, 1, 10],
                               'polynomial__degree': [1, 2, 3, 4, 5]})
```

```
[41]: # find the best model:
      best=search.best_estimator_
      best
```

```
[41]: Pipeline(steps=[('polynomial',
                       PolynomialFeatures(degree=3, include_bias=False)),
                      ('ss', StandardScaler()), ('model', Lasso(alpha=10, tol=0.2))])
```

```
[42]: # calculate the   2 on the test data
      best.score(X_test,y_test)
```

```
[42]: 0.9390483610147886
```

Elastic Net

In this section, let's review the Elastic Net Regression. It combines L1 and L2 priors as regularizes or penalties. So, we can combine the two as follows:

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} ||\mathbf{y} - \mathbf{Xw}||^2\_2 + \alpha\rho||\mathbf{w}||\_10.5\alpha(1-\rho)||\mathbf{w}||^2\_2$$

Additionally to the alpha term ($\alpha$), we have a mixing parameter, $\rho$, such that $0 \le \rho \le 1$. For $\rho$=0, the penalty is an L2 regularization . For $\rho = 0$, it is L1 regularization; otherwise, it is a combination of L1 and L2. In *scikit-learn* the parameter is called `l1_ratio`. Unlike the Ridge Regression, Elastic Net finds zero coefficients. In many cases Elastic Net performs better than Lasso, as it includes features that are correlated with one another. One drawback of the Elastic Net is you have two hyperparameters. Lets create a model where `alpha=0.1`and `l1_ratio=0.5` and fit the data with this model.

```
[43]: enet = ElasticNet(alpha=0.1, l1_ratio=0.5)
      enet.fit(X_train,y_train)
```

```
[43]: ElasticNet(alpha=0.1)
```

```
[44]: predicted=enet.predict(X_test)
      predicted
```

```
[44]: array([27876.23069332, 19692.47948228, 11278.55904242, 11378.11762817,
             25990.32173225,  5511.46549415,  7591.3310787 ,  7796.50950874,
              9762.62529929,  8646.72495203, 17018.24269379,  7412.95036812,
             17031.98094445, 10732.9267982 , 37491.12185882,  5254.66797235,
              -286.43084016, 16075.95423556, 11330.87482645, 10579.46380101,
             11539.66267619])
```

```
[45]: # calculate the  2 on the test data
      print("R^2 on training  data ", enet.score(X_train, y_train))
      print("R^2 on testing data ", enet.score(X_test,y_test))
```

```
R^2 on training  data  0.8726983414262406
R^2 on testing data  0.9289334382162673
```

If we compare the Elastic Net to Lasso Regression and Ridge Regression, we see the results on the $R^2$ are better than the Elastic Net and many of the coefficients are zero.

```
[46]: plot_coef(X,la,name="Lasso Regression")
      plot_coef(X,enet,name="Elastic net ")
```

## Lasso Regression



R^2 on training  data  0.799083710857619
R^2 on testing data  0.910818845927178

Elastic net

```
R^2 on training  data  0.8726983414262406
R^2 on testing data  0.9289334382162673
```

Exercise 3

Create and fit the Elastic Net model and the Ridge Regression models and plot the coefficients for both models using the `plot_coef()` function.

```
[47]: enet = ElasticNet(alpha=0.01, l1_ratio=0)
      enet.fit(X_train,y_train)
      rr = Ridge(alpha=0.01)
      rr.fit(X_train,y_train)
      plot_coef(X,rr,name="Ridge Regression")

      plot_coef(X,enet,name="Elastic net l1_ratio=0 ")
```

Ridge Regression

R^2 on training  data  0.9091956531801182
R^2 on testing data  0.9478784615596494

Elastic net l1_ratio=0

```
R^2 on training  data  0.8941741868792897
R^2 on testing data  0.940724027816336
```

Exercise 4

Create a Pipeline object, apply polynomial features (degree = 2), perform data standardization, then apply Elastic Net with `alpha=0.1` and `l1_ratio=0.1` parameters. Fit the model using the training data, then calculate the $R^2$ on the training and testing data.

```
[48]: Input=[ ('polynomial',␣
      ↪PolynomialFeatures(include_bias=False,degree=2)),('ss',StandardScaler() ),␣
      ↪('model',ElasticNet(alpha=0.1, l1_ratio=0.1))]
      pipe = Pipeline(Input)
      pipe.fit(X_train, y_train)
      print("R^2 on training  data ",pipe.score(X_train, y_train))
      print("R^2 on testing data ",pipe.score(X_test,y_test))
```

```
R^2 on training  data  0.9703466188354504
R^2 on testing data  0.9474289882341259
```

Exercise 5

Search for the best combination of hyperparameters by creating a `GridSearchCV` object for Elastic

Net Regression. Find the best parameter values using the pipeline object, as used in the above examples. Use `param_grid`, then find thee $R^2$ on the test data using the best estimator.

```
[49]: param_grid = {
          "polynomial__degree": [ 1, 2,3,4,5],
          "model__alpha":[0.0001,0.001,0.01,0.1,1,10],
          "model__l1_ratio":[0.1,0.25,0.5,0.75,0.9]
      }
```

```
[50]: Input=[ ('polynomial',␣
      ↪PolynomialFeatures(include_bias=False,degree=2)),('ss',StandardScaler() ),␣
      ↪('model',ElasticNet(tol = 0.2))]
      pipe = Pipeline(Input)
      search = GridSearchCV(pipe, param_grid, n_jobs=2)
      search.fit(X_test, y_test)
      best=search.best_estimator_
      best.score(X_test,y_test)
```

```
[50]: 0.9800717238596847
```

---

Principal Component Analysis (PCA)

In this example, we will explore Principal Component Analysis to reduce the dimensionality of our data. We will do so by creating a Pipeline object first, then applying standard scaling and performing PCA, and then applying Elastic Net Regularization with the following parametrs: `tol=0.2`, `alpha=0.1` and `l1_ratio=0.1`. Finally, we will fit the model using the training data, then calculate the $R^2$ on the training and testing data sets.

Before adding PCA as a prep-processing step, we have to standardize our data. Scaling the features makes them have the same standard deviation.

```
[51]: scaler = StandardScaler()
      X_train[:] = scaler.fit_transform(X_train)
      X_train.columns = [f'{c} (scaled)' for c in X_train.columns]
```

```
[52]: # perform PCA
      pca = PCA()
      pca.fit(X_train)
```

```
[52]: PCA()
```

We can find the projection of the dataset onto the principal components, let's call it X_train_hat , this is our "new" dataset, it is the same shape as the original dataset.

```
[53]: X_train_hat = pca.transform(X_train)
      print(X_train_hat.shape)
```

```
(184, 35)
```

```
[54]: # make the new dataset to a dataframe
      X_train_hat_PCA = pd.DataFrame(columns=[f'Projection  on Component {i+1}' for i
        ↪in range(len(X_train.columns))], data=X_train_hat)
      X_train_hat_PCA.head()
```

```
[54]:    Projection  on Component 1  Projection  on Component 2  \
      0                    3.294148                   -1.989679
      1                   -2.801332                   -0.496152
      2                   -5.993408                   -0.535847
      3                   -3.462434                    0.060153
      4                    4.664168                   -0.948145

         Projection  on Component 3  Projection  on Component 4  \
      0                   -2.182446                   -1.156530
      1                   -0.654240                    0.175846
      2                   -0.270387                   -0.833262
      3                   -0.841084                   -0.505797
      4                   -1.440491                   -0.739345

         Projection  on Component 5  Projection  on Component 6  \
      0                   -0.715546                   -0.605020
      1                    0.267657                    0.914558
      2                    0.898968                    1.739295
      3                   -0.771184                    0.138637
      4                   -1.235190                   -0.889695

         Projection  on Component 7  Projection  on Component 8  \
      0                   -1.628886                    0.853635
      1                    0.066678                    1.120531
      2                   -0.624702                    0.171002
      3                   -0.214593                    0.253033
      4                   -0.726741                   -0.228066

         Projection  on Component 9  Projection  on Component 10  …  \
      0                    2.355868                     0.434029  …
      1                   -1.804675                    -2.048843  …
      2                    0.741584                    -0.811344  …
      3                   -0.173882                     0.780597  …
      4                   -0.883634                    -0.048221  …

         Projection  on Component 26  Projection  on Component 27  \
      0                    0.116948                     0.089306
      1                    0.221664                    -0.027223
      2                    0.014135                     0.207760
      3                    0.026002                    -0.146690
      4                   -0.005600                     0.124976
```

```
      Projection  on Component 28  Projection  on Component 29  \
0                   1.051890e-15                  -7.248345e-16
1                   8.320969e-16                  -6.152919e-16
2                  -1.564247e-15                  -1.113952e-15
3                   1.670410e-16                   8.350073e-17
4                  -8.961753e-16                   4.920165e-16

      Projection  on Component 30  Projection  on Component 31  \
0                  -2.655793e-16                   1.362665e-16
1                   4.550324e-16                  -9.757210e-16
2                  -4.280176e-16                   9.804134e-17
3                   2.622225e-16                  -2.194562e-16
4                  -7.972540e-16                   7.799564e-16

      Projection  on Component 32  Projection  on Component 33  \
0                  -1.238006e-16                  -1.465801e-15
1                  -4.277862e-16                  -5.731525e-16
2                   1.767552e-16                  -5.172664e-16
3                   4.644181e-17                  -8.653447e-16
4                   2.378532e-16                   1.035291e-16

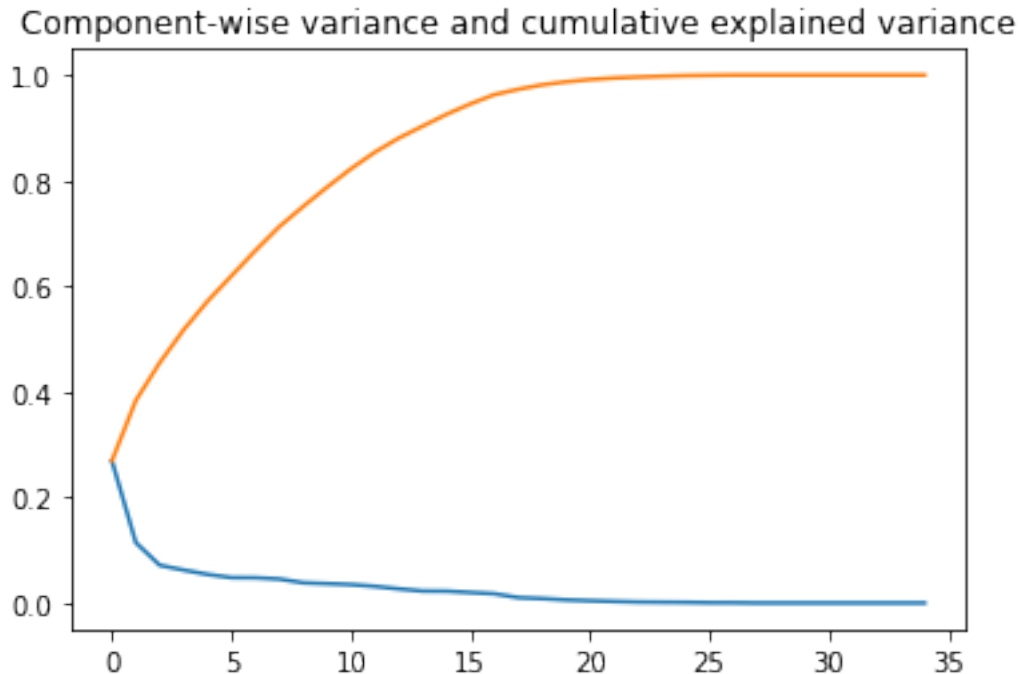      Projection  on Component 34  Projection  on Component 35
0                   2.820806e-16                   1.659746e-16
1                  -5.517287e-16                  -4.646918e-16
2                  -7.738870e-16                   2.370532e-16
3                  -4.613195e-16                   1.403707e-17
4                   2.139927e-16                   4.569918e-16

[5 rows x 35 columns]
```

Now, let's see how much variance can be explained using these principal components (PCs);

```
[55]: plt.plot(pca.explained_variance_ratio_)
      plt.plot(np.cumsum(pca.explained_variance_ratio_))
      plt.title("Component-wise variance and cumulative explained variance")
```

```
[55]: Text(0.5, 1.0, 'Component-wise variance and cumulative explained variance')
```

Component-wise variance and cumulative explained variance

In the graph above, the component-wise variance is depicted by the blue line, and the cumulative explained variance is explained by the orange line. We are able to explain ~100% of the variance using just the first 20 PCs. Let's filter our dataset down to these 20 PCs.

```
[56]: N = 20
      X_train_hat_PCA = X_train_hat_PCA.iloc[:, :N]
```

Let's create an Elastic Net model where `tol=0.2`, `alpha=0.1` and `l1_ratio=0.1` and fit the data with this model.

```
[57]: enet = ElasticNet(tol = 0.2, alpha=100, l1_ratio=0.75)
      enet.fit(X_train_hat_PCA, y_train)
```

```
[57]: ElasticNet(alpha=100, l1_ratio=0.75, tol=0.2)
```

Exercise 6

In this Exercise, create a Pipeline object, apply standard scaling, perform PCA and then finally fit an Elastic Net with `tol=0.2`, `alpha=0.1` and `l1_ratio=0.1` parameters. Calculate the scores, $R^2$, on the training and testing data sets.

```
[58]: Input = [
          ('scaler', StandardScaler()),
          ('pca', PCA(n_components = N)),
          ('model', ElasticNet(tol =0.2, alpha=0.1, l1_ratio=0.1))
              ]
```

```
pipe = Pipeline(Input)
pipe.fit(X_train, y_train)
print("R^2 on training  data ", pipe.score(X_train, y_train))
print("R^2 on testing data ", pipe.score(X_test,y_test))
```

```
R^2 on training  data  0.8829912421979917
R^2 on testing data  -74934.50860233128
```

**Conclusion on the above**: The R-squared ($R^2$) value on the training data is 0.883, which indicates that the model explains 88.3% of the variance in the training data. However, the $R^2$ value on the testing data is -74934.51, which indicates that the model performs very poorly on the testing data and is likely overfitting to the training data. This means that the model may not generalize well to new, unseen data.