

## Ridge Regression

When the numerical features of the data in question are collinear, there arises a few issues of numerical instability and generalization. Ordinary least squares does not handle this well and in order to mitigate this we add a regularization term that penalizes data points that may represent outliers in the dataset

```
In [1]: import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
from sklearn.datasets import make_regression

In [2]: #This regression method is more apt for when a certain data point is an anomaly
and may skew the results
#we add a regularization paramater that helps control that
import numpy as np
def RidgeRegression(A,b,lambd_):
    #lambda_ here is the reguralization parameter,also solving Ax=b
    n, m = A.shape
    I = np.identity(m)
    x= np.dot(np.dot(np.linalg.inv(np.dot(A.T, A) + lambd_ * I), A.T), b)
    return x
```

Creating the data set below

```
In [3]: X, y, coefficients = make_regression(
    n_samples=50,
    n_features=1,
    n_informative=1,
    n_targets=1,
    noise=5,
    coef=True,
    random_state=1
)
```

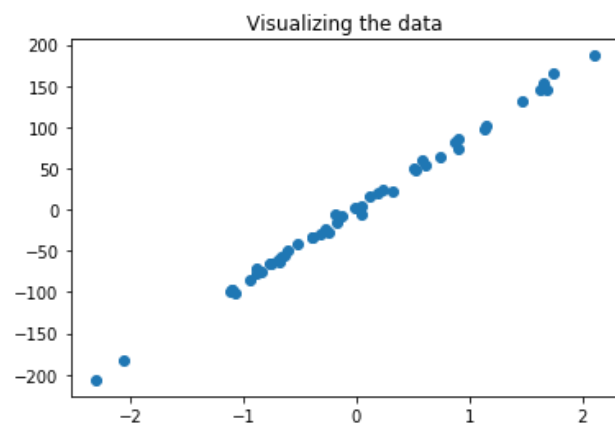
## Note on Mean sqared error (MSE)

The MSE is a measure of the quality of an estimator. It calculates the mean squared difference across two quantities. The mean squared error is never negative, and the closer to zero, the better the estimate

```
In [4]: def MSE(y_true,y_predicted):
    return (1/len(y_true))*np.sum((np.subtract(y_true,y_predicted))**2)
```

Let us first visualize the data, the noise is not very high in this sample

```
In [5]: plt.scatter(X, y)
plt.title("Visualizing the data")
plt.show()
```



Now let us plot our estimate using various penalty coefficients

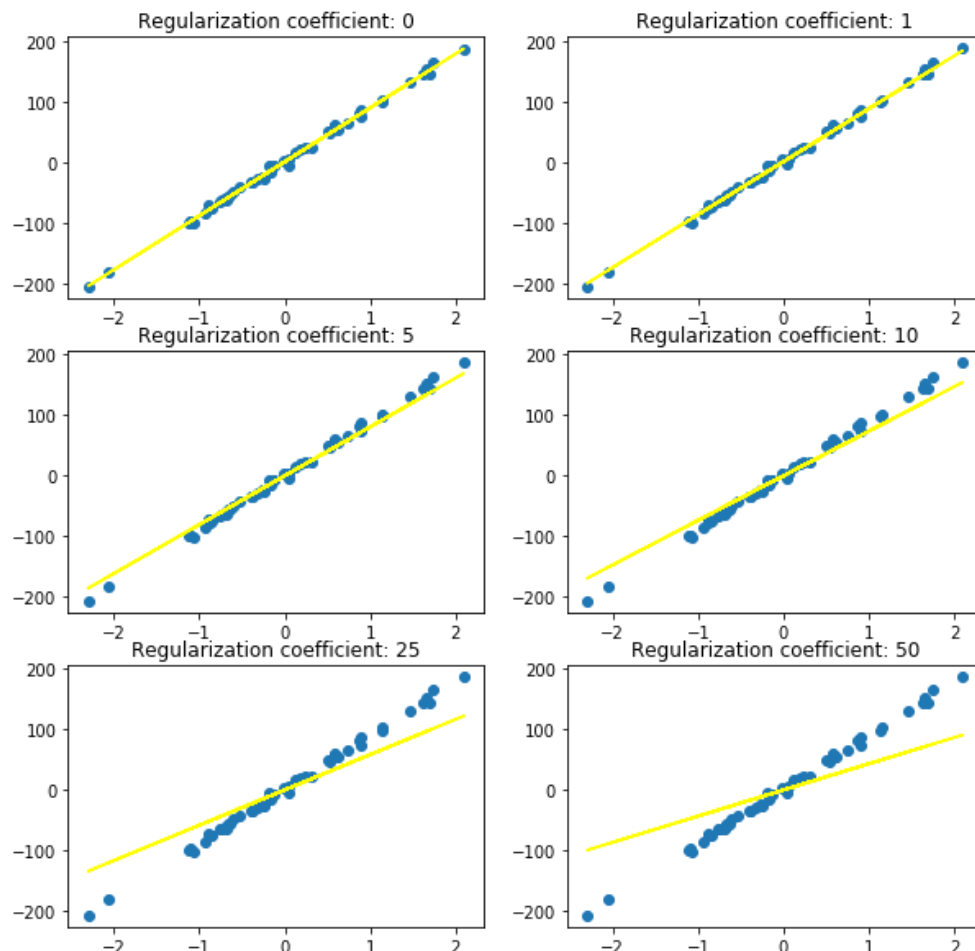
```

In [6]: lambdas= [0,1,5,10,25,50]
plt.figure(figsize=(10,10))

for i in range(len(lambdas)):
    w= RidgeRegression(X,y,lambdas[i])
    pred_y= np.matmul(X,w)
    plt.subplot(3,2,i+1)
    plt.title("Regularization coefficient: " + str(lambdas[i]))
    plt.scatter(X, y)
    plt.plot(X, pred_y, c='yellow')
    mse= MSE(y,pred_y)
    print("The mean squared error using the penalty coefficient " + str(lambdas
[i]) + " is " + str(mse) )
plt.show()

```

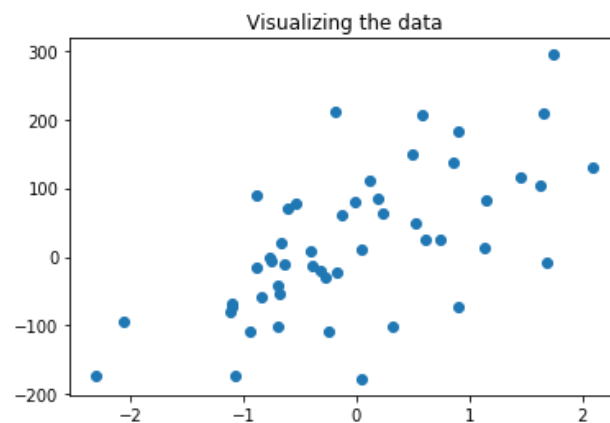
The mean squared error using the penalty coefficient 0 is 17.779274002775864  
 The mean squared error using the penalty coefficient 1 is 21.025088343365702  
 The mean squared error using the penalty coefficient 5 is 86.92923432921926  
 The mean squared error using the penalty coefficient 10 is 248.0104499012563  
 The mean squared error using the penalty coefficient 25 is 919.8665198890254  
 The mean squared error using the penalty coefficient 50 is 2006.3785724573142



So what's going on here, since the noise in the data set isn't very high, the regularization coefficient seems to be doing very little as it increases, at it unnecessarily penalizes points that shouldn't be penalized, but now, let us really increase the noise in the data set through outliers and observe what occurs.

```
In [7]: x, Y, coefficients = make_regression(
        n_samples=50,
        n_features=1,
        n_informative=1,
        n_targets=1,
        noise=100,
        coef=True,
        random_state=1
    )
```

```
In [8]: plt.scatter(x, Y)
        plt.title("Visualizing the data")
        plt.show()
```



As we can now see the data set is much noisier, let us see how our regularization parameters perform now.

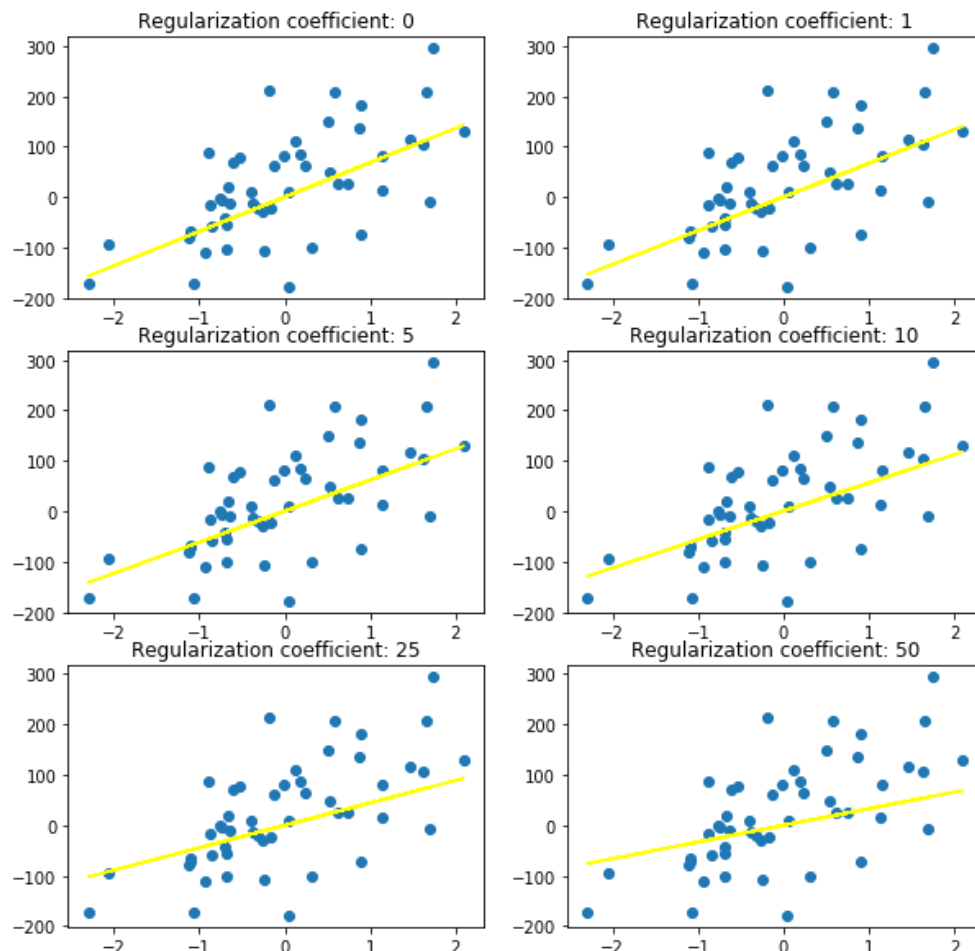
```

In [9]: lambdas= [0,1,5,10,25,50]
plt.figure(figsize=(10,10))

for i in range(len(lambdas)):
    w= RidgeRegression(x,Y,lambdas[i])
    pred_y= np.matmul(x,w)
    plt.subplot(3,2,i+1)
    plt.title("Regularization coefficient: " + str(lambdas[i]))
    plt.scatter(x, Y)
    plt.plot(x, pred_y, c='yellow')
    mse= MSE(Y,pred_y)
    print("The mean squared error using the penalty coefficient " + str(lambdas
[i]) + " is " + str(mse) )
plt.show()

```

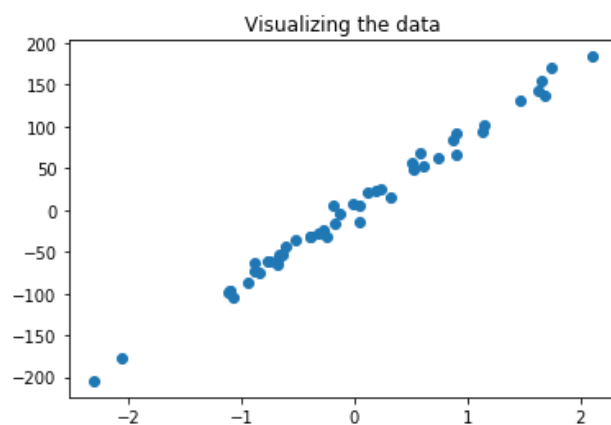
The mean squared error using the penalty coefficient 0 is 7111.709601110346  
 The mean squared error using the penalty coefficient 1 is 7113.601151011548  
 The mean squared error using the penalty coefficient 5 is 7152.007839671917  
 The mean squared error using the penalty coefficient 10 is 7245.880478531952  
 The mean squared error using the penalty coefficient 25 is 7637.415282884596  
 The mean squared error using the penalty coefficient 50 is 8270.59745229904



As we can know see the regularization parameters offer similar performance in the case of increased noise, however, a higher regularization parameter is not necessarily better, let us look at one last data set, with relatively homogenous data, with a few outliers

```
In [10]: X, Y, coefficients = make_regression(
          n_samples=50,
          n_features=1,
          n_informative=1,
          n_targets=1,
          noise=10,
          coef=True,
          random_state=1
        )
```

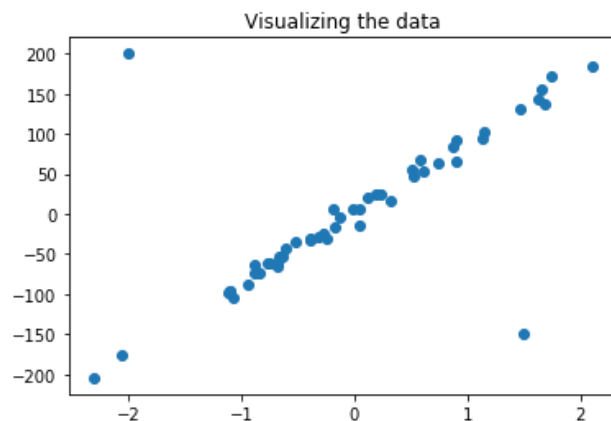
```
In [11]: plt.scatter(X, Y)
          plt.title("Visualizing the data")
          plt.show()
```



```
In [12]: Y= np.append(Y, [200, -150])
          X= np.append(X, [-2, 1.5])
          X= np.expand_dims(X, axis=1)
          print(np.shape(X))
          print(np.shape(Y))
```

```
(52, 1)
(52,)
```

```
In [13]: plt.scatter(X, Y)
plt.title("Visualizing the data")
plt.show()
```



So now, we have some definitive outliers, let us see what our regularization coefficients do now:

```
In [14]: lambdas= [0,1,5,10,25,50]
plt.figure(figsize=(10,10))

for i in range(len(lambdas)):
    w= RidgeRegression(X,Y,lambdas[i])
    pred_y= np.matmul(X,w)
    plt.subplot(3,2,i+1)
    plt.title("Regularization coefficient: " + str(lambdas[i]))
    plt.scatter(X, Y)
    plt.plot(X, pred_y, c='yellow')
    mse= MSE(Y,pred_y)
    print("The mean squared error using the penalty coefficient " + str(lambdas
[i]) + " is " + str(mse) )

plt.show()
```

The mean squared error using the penalty coefficient 0 is 3822.9132381893432  
The mean squared error using the penalty coefficient 1 is 3824.4303500563256  
The mean squared error using the penalty coefficient 5 is 3855.8140551973756  
The mean squared error using the penalty coefficient 10 is 3934.5434909947685  
The mean squared error using the penalty coefficient 25 is 4278.859367887072  
The mean squared error using the penalty coefficient 50 is 4870.677139425841

