# Week 3 Assignment

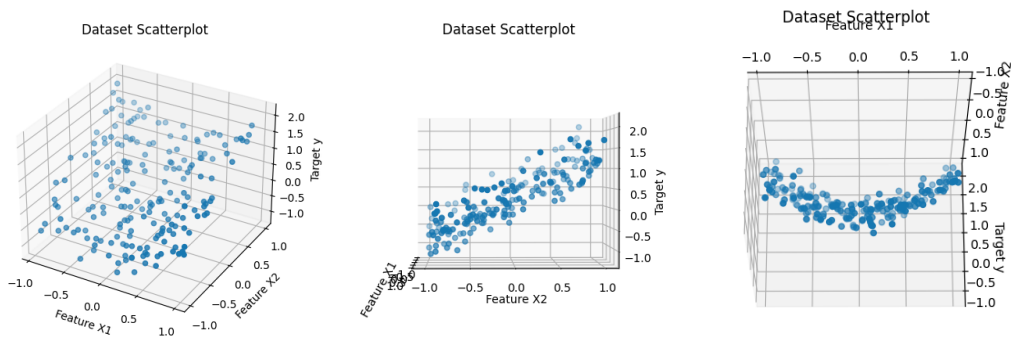CS7CS4/CSU44061 Machine Learning

Cormac Madden – 18319983

Dataset ID: # id:11-11-11

1)
   a) I plotted the data using the code provided in the assignment.

   It appears the training data lies on a curved plane. From the X2,Y perspective the data appears linear (Fig 2), however from another perspective the data appears to follow a quadratic curve (Fig 3).



   b) Below are the features of the model for the value of C:1,10,100,1000,10000



   $J(\theta)$ = mean squared error

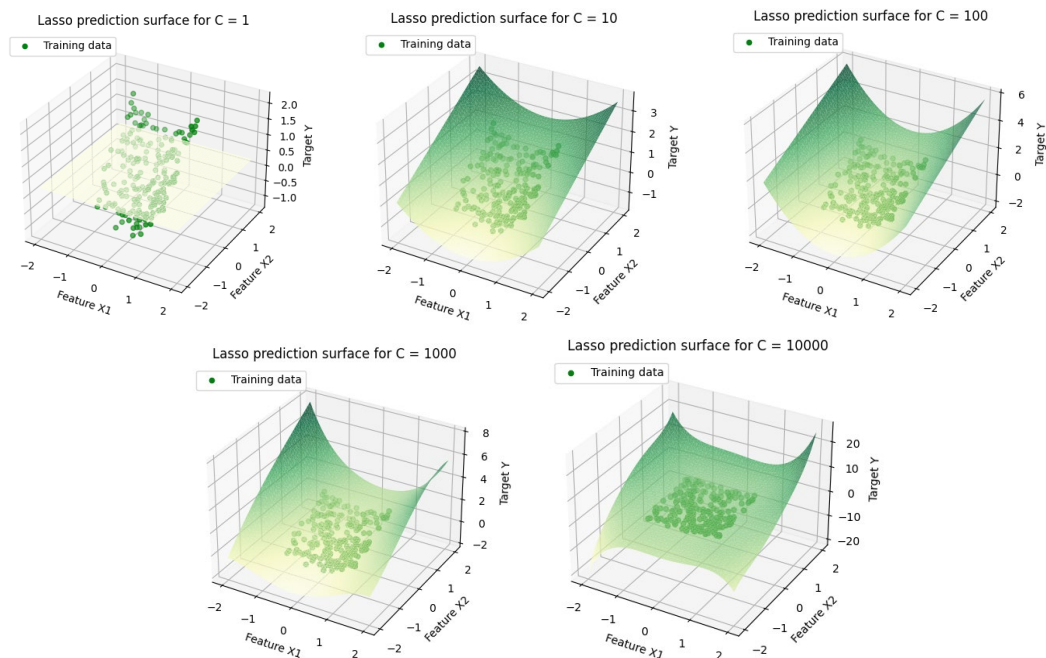| C | 1 | 10 | 100 | 1000 | 10000 |
|---|---|----|-----|------|-------|
| $J(\theta)$ | 0.488932 | 0.071785 | 0.037024 | 0.036140 | 0.035303 |

The first output has 0 features because the penalty of C=1 is too strong. It gives the same mean squared error as the baseline (0.488932) which is just the mean. When C=1 there are two features, for b and $a^2$. This reduces the error significantly to 0.071785, and when drawn on a graph would show the curve we predicted in a).

For C=100-10000 more features are added, which adds more curves to the prediction model. This reduces the error but may also be unnecessary and be over-fitting the data. A simpler model such as for C=100 is probably adequate, as the decrease in error is marginal after that.

Below is the code used to create the data above:

```
1   #Xpoly
2   Xpoly = PolynomialFeatures(degree = 5).fit_transform(X)
3
4   baseline = DummyRegressor(strategy="mean").fit(Xpoly, y)
5   print("J(0_baseline) = %f\n"%mean_squared_error(y, baseline.predict(Xpoly)))
6
7   C = [1,10,100,1000,10000]
8   lassos = []
9   for Ci in C:
10      model = Lasso(alpha=1/(2*Ci)).fit(Xpoly,y)
11      lassos.append(model)
12      print("\nC value: "+str(Ci))
13      print("Lasso coef: "+str(model.coef_))
14      print("Lasso intercept: "+str(model.intercept_))
15      print("J(0) = %f\n"%mean_squared_error(y, model.predict(Xpoly)))
16
```

c)



I decided to show the predictions on several different plots, in order to make it easier to read.

As C is increased the curve fits more tightly to the training data, but as the plane extends from the training data is bends sharply in unexpected ways. Note the difference in scale on the Target Y axis on the five plots. It is possible but highly unlikely that this is representative of the data. Where C=10 or 100 it looks as we had anticipated in question a) and would extend out in that logical direction.
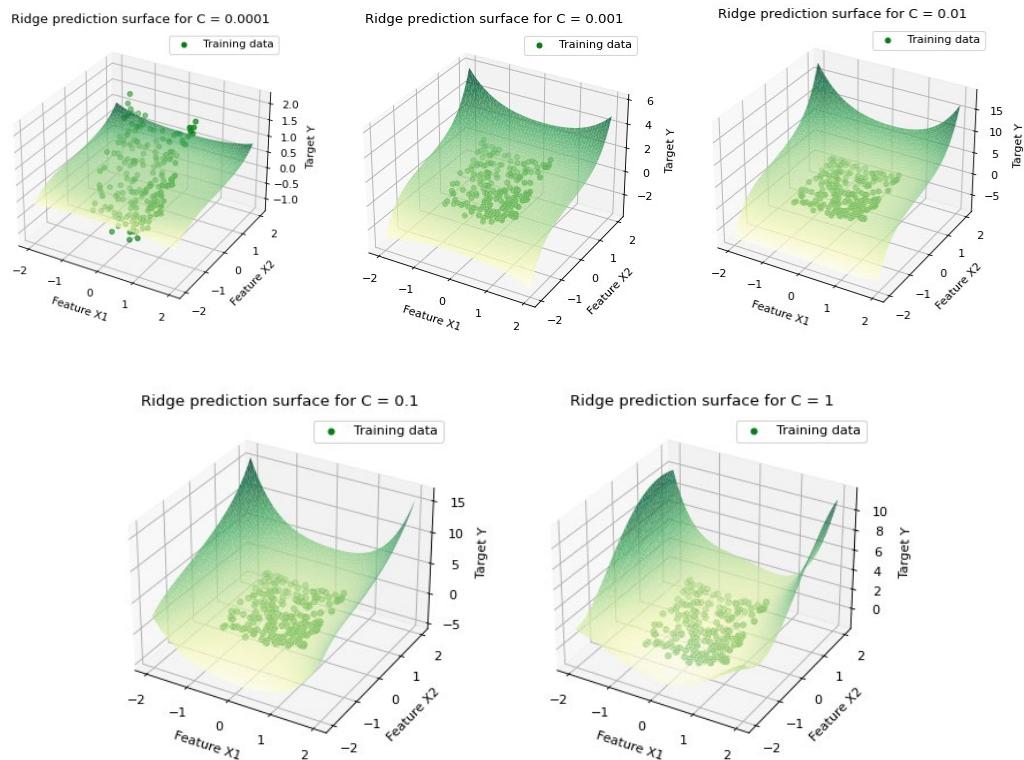
For the higher C values, smaller clusters of data at the edge of the dataset are having a significant impact on the direction of the extremities of the prediction plane.

d) Under-fitting is when the model is too simple to represent the data. It will struggle to predict the data and example of this can be seen in the first graph from question c).

Over-fitting is where the model tries to adapt to too much of the training data and as a result performs poorly on new data. An example of this is in the last graph from question c)

The penalty weight parameter C can be used to manage how closely the model fits the training data. As can be seen in the figures above, when a small value of C is used, the model uses less features and can tend towards underfitting the data. Whereas when large values of C are used the model uses more features and tends towards overfitting the data.

e) The following figures are the graphs produced from creating a Ridge Regression model and using various values for C. C in this case affects the cost function and is referred to as a hyperparameter.

For this model all the features are used but the penalty affects the impact the training data will have on the prediction curve. As you can see from the data all the features in the first model are close to zero, whereas in the later models some of the feaures increase quite significantly.



The mean square error J(θ) for each Ridge model was also found using the same baseline as in question (i)(b) and (c), and is shown in the table below.

| C | 0.0001 | 0.001 | 0.01 | 0.1 | 1 |
|---|---|---|---|---|---|
| J(θ) | 0.468832 | 0.341606 | 0.106170 | 0.042716 | 0.036295 |

There is a notable difference between the Ridge and Lasso model, in that the Lasso model will cancel small parameters which makes it easy to classify the data as linear, quadratic etc. The Ridge model on the other hand only adjusts the significance of each parameter without cancelling any. This means the model will always be more complex and difficult to interpret and manually manipulate.

2)
   a) Below is a plot of the mean and the standard deviation of the prediction error, on the Lasso Regression model, against various values for C.
      The range of values for C were chosen by repeatedly creating and plotting the model to find a reasonable range of values that would illustrate a clear trend.



Lasso Regression Cross-Validation of C values

   b) Based off the graph and the data presented earlier I would recommend using a C value of 100. We know from earlier that a C of 100 only uses three features so it is still a simple model, and the mean error and standard deviation is adequately low.

c) Below is a plot of the mean and the standard deviation of the prediction error, on the Ridge Regression model, against various values for C.

The range of values for C were chosen by repeatedly creating and plotting the model to find a reasonable range of values that would illustrate a clear trend.



Based off the graph and the data presented earlier I would recommend using a C value of 0.1. It has a small standard deviation and appears to be a reasonable balance between representing the data and overfitting.

Appendix

```python
# %%
%matplotlib widget
import numpy as np
import math
from matplotlib import cm
import pandas as pd
import scipy as sp
import seaborn as sn
import matplotlib.pyplot as plt
import matplotlib as matlib
from sklearn.linear_model import LinearRegression
from sklearn.linear_model import LogisticRegression
from sklearn.linear_model import Ridge

from sklearn.model_selection import KFold
from sklearn.model_selection import RepeatedKFold
from sklearn.model_selection import train_test_split

from sklearn.svm import LinearSVC
from sklearn.dummy import DummyRegressor
from sklearn.metrics import mean_squared_error
from sklearn.linear_model import Lasso
from sklearn.preprocessing import PolynomialFeatures
from mpl_toolkits.mplot3d import Axes3D


df = pd.read_csv('week3.csv')
df.columns = ["X1","X2","y"]
X1=df.iloc[:,0]
X2=df.iloc[:,1]
y=df.iloc[:,2]
X=np.column_stack((X1,X2))

# %% [markdown]
# a i)

# %%
fig = plt.figure()
ax = fig.add_subplot(111, projection='3d')
ax.scatter(X1,X2,y)
ax.set_xlabel('Feature X1')
ax.set_ylabel('Feature X2')
ax.set_zlabel('Target y')
plt.title('Dataset Scatterplot')
plt.show()

# %% [markdown]
# a) ii)

# %%
Xpoly = PolynomialFeatures(degree = 5).fit_transform(X)

baseline = DummyRegressor(strategy="mean").fit(Xpoly, y)
print("J(θ_baseline) = %f\n"%mean_squared_error(y, baseline.predict(Xpoly)))


C = [1,10,100,1000,10000]
lassos = []
for Ci in C:
    model = Lasso(alpha=1/(2*Ci)).fit(Xpoly,y)
    lassos.append(model)
    print("\nC value: "+str(Ci))
    print("Lasso coef: "+str(model.coef_))
```
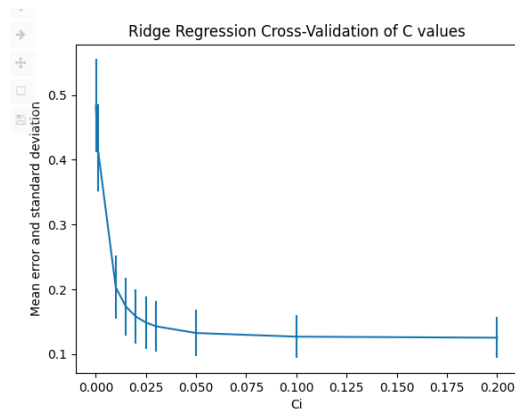
```python
    print("Lasso intercept: "+str(model.intercept_))
    print("J(θ) = %f\n"%mean_squared_error(y, model.predict(Xpoly)))


# %%
### i) c)
Xtest = []
grid = np.linspace(-2,2)
for i in grid:
    for j in grid:
        Xtest.append([i,j])

Xtest = np.array(Xtest)
Xtest = PolynomialFeatures(5).fit_transform(Xtest)

from matplotlib import cm
Xpoly = PolynomialFeatures(5).fit_transform(X)

C_range = [1, 10, 100, 1000, 10000]
for Ci in C_range:
    model = Lasso(alpha=1/(2*Ci))
    model.fit(Xpoly, y)
    y_pred=model.predict(Xtest)

    fig = plt.figure(dpi=100)
    ax = fig.add_subplot(111, projection='3d')

    ax.scatter(X1, X2, y, c='g', label="Training data")
    surf = ax.plot_trisurf(Xtest[:,1], Xtest[:,2], y_pred, cmap=cm.YlGn, alpha=0.8, linewidth=0,
antialiased=True)
    ax.set_title('Lasso prediction surface for C = %.0f'%Ci)
    ax.set(xlabel='X1', ylabel='X2', zlabel='Target Y')
    ax.legend(loc='upper left')

    plt.xlabel("Feature X1"); plt.ylabel("Feature X2")

    plt.show()


# %%
# (i)(e)

baseline = DummyRegressor(strategy="mean").fit(Xpoly, y)
print("J(θ_baseline) = %f\n"%mean_squared_error(y, baseline.predict(Xpoly)))

ridges = []
C_ridge = [0.0001,0.001,0.01,0.1,1,10]
for Ci in C_ridge:
    model = Ridge(alpha=1/(2*Ci)).fit(Xpoly, y)
    theta = np.insert(model.coef_, 0, model.intercept_)
    ridges.append(model)

    print("\nC value: "+str(Ci))
    print("θ =", theta)
    print("J(θ) = %f\n"%mean_squared_error(y, model.predict(Xpoly)))

    fig = plt.figure(dpi=80)
    ax = fig.add_subplot(111, projection='3d')

    ax.scatter(X[:,0], X[:,1], y, c='g', label="Training data")
    surf = ax.plot_trisurf(Xtest[:,1], Xtest[:,2], model.predict(Xtest), cmap=cm.YlGn, alpha=0.8,
linewidth=0, antialiased=True)

    ax.set_title('Ridge prediction surface for C = ' + str(Ci))
```

```python
    ax.set(xlabel='Feature X1', ylabel='Feature X2', zlabel='Target Y')
    ax.legend()

    plt.show()

print("exit")


# %%
### ii a)

kf = KFold(n_splits=5)
mean_error=[]
std_error=[]
Cs2 = [1,2,3,4,5,6,7,8,9,10,20,30,40,50,100]
for Ci in Cs2:
    model = Lasso(alpha=1/(2*Ci))
    print("\nC: "+str(Ci))
    kf = KFold(n_splits=5)
    temp = []
    for train, test in kf.split(X):
        model.fit(X[train],y[train])
        ypred = model.predict(X[test])
        print("Lasso coeff: "+str(model.coef_))

        from sklearn.metrics import mean_squared_error
        temp.append(mean_squared_error(y[test],ypred))
    mean_error.append(np.array(temp).mean())
    std_error.append(np.array(temp).std())

plt.figure(dpi=100)
plt.errorbar(Cs2,mean_error,yerr=std_error)
plt.xlabel('Ci'); plt.ylabel('Mean error and standard deviation')
plt.title('Lasso Regression Cross-Validation of C values')
plt.show()

### ii c)
mean_error=[]
std_error=[]
C_ridge2 = [0.0001,0.001,0.01,0.015,0.02,0.025,0.03,0.05,0.1,0.2]
for Ci in C_ridge2:
    model = Ridge(alpha=1/(2*Ci))

    from sklearn.model_selection import KFold
    kf = KFold(n_splits=5)
    temp = []
    print("\nC: "+str(Ci))

    for train, test in kf.split(X):
        model.fit(X[train],y[train])
        ypred = model.predict(X[test])
        print("Ridge coeff: "+str(model.coef_))
        from sklearn.metrics import mean_squared_error
        temp.append(mean_squared_error(y[test],ypred))
    mean_error.append(np.array(temp).mean())
    std_error.append(np.array(temp).std())

plt.figure(dpi=100)
plt.errorbar(C_ridge2,mean_error,yerr=std_error)
plt.xlabel('Ci'); plt.ylabel('Mean error and standard deviation')
plt.title('Ridge Regression Cross-Validation of C values')
plt.show()
```