

# STA 208: Homework 1 (Do not distribute)

Due 04/24/2022 midnight (11:59pm)

## Instructions:

1. Submit your homework using one file name "Lastname\_Firstname\_hw1.html" on canvas.
2. The written portions can be either done in markdown and TeX in new cells or written by hand and scanned. Using TeX is strongly preferred. However, if you have scanned solutions for handwriting, you can submit a zip file. Please make sure your handwriting is clear and readable and your scanned files are displayed properly in your jury notebook.
3. Your code should be readable: writing a piece of code should be compared to writing a page of a book. Adopt the one-statement-per-line rule. Consider splitting a lengthy statement into multiple lines to improve readability. (You will lose one point for each line that does not follow the one-statement-per-line rule)
4. To help understand and maintain code, you should always add comments to explain your code. (homework with no comments will receive 0 points). For a very long comment, please break it into multiple lines.
5. In your Jupyter Notebook, put your answers in new cells after each exercise. You can make as many new cells as you like. Use code cells for code and Markdown cells for text.
6. Please make sure to print out the necessary results to avoid losing points. We should not run your code to figure out your answers.
7. However, also make sure we are able to open this notebook and run everything here by running the cells in sequence; in case that the TA wants to check the details.
8. You will be graded on correctness of your math, code efficiency and succinctness, and conclusions and modelling decisions

## Exercise 1 (Empirical risk minimization) (20 pts, 5 pts each)

Consider Poisson model with rate parameter  $\lambda$  which has PMF,

$$p(y|\lambda) = \frac{\lambda^y}{y!} e^{-\lambda}$$

where  $y = 0, 1, \dots$  is some count variable. In Poisson regression, we model  $\lambda = e^{\beta^T x}$  to obtain  $p(y|x, \beta)$ .

1. Let the loss function for Poisson regression be  $\ell(\beta) \propto -\log \phi(y_i|x_i, \beta)$  for a dataset consisting of predictor variables and count values  $\{x_i, y_i\}_{i=1}^n$ . Here  $\propto$  means that we disregard any additive terms that are not dependent on  $\beta$ . Write an expression for  $\ell_i$  and derive its gradient.

**Solution:**

Based on the above expression, we can find the log density of the poisson distribution after substituting  $\lambda = e^{\beta^T x}$  back into the probability function. On doing so, we get the density function as:

$$p(y|x, \beta) = \frac{(e^{\beta^T x})^y}{y!} e^{-e^{\beta^T x}}$$

On taking the log of this, we get:

$$\log p(y|x, \beta) = y\beta^T x - \log(y!) - e^{\beta^T x}$$

Based on the definition of the proportionality, we just ignore the additive terms not dependent on  $\beta$  to get the loss function which is given as:

$$J(\beta) = -\log \phi(y_i|x_i, \beta) = -(\gamma \beta^T x_i - e^{\beta^T x_i}) = e^{\beta^T x_i} - y \beta^T x_i$$

The gradient of this loss function is given as:

$$\frac{\partial J(\beta)}{\partial \beta} = x_i(e^{\beta^T x_i} - y_i)$$

1. Show that the empirical risk  $J_n(\beta)$  is a convex function of  $\beta$ .

**Solution:**

The risk function  $J_n(\beta)$  is defined as summation of loss functions for each data point. It is formally defined as:

$$J_n(\beta) = \sum_{i=1}^n J(\beta)$$

In order to show that the risk function is convex, we can use the fact that if loss function is convex, then the risk function is convex as risk function is a linear combination or sum of the loss functions at every  $i$ .

Now in order to show the loss function is convex, we have to show that the double derivative of the loss function is  $\geq 0$ .

From part 1, we have the gradient as:

$$\frac{\partial J(\beta)}{\partial \beta} = x_i(e^{\beta^T x_i} - y_i)$$

Taking the double derivate of the above expression as:

$$\frac{\partial^2 J(\beta)}{\partial \beta^2} = x_i^T x_i e^{\beta^T x_i} \geq 0$$

Since for any values of  $x_i$ , the double derivative would be greater than 0, we can conclude that the loss function is convex.

Since the loss function is convex and the risk function is the sum of loss functions, we can conclude that the risk function is convex as well.

1. Consider the mapping  $F_n(\beta) = \beta - \eta \nabla J_n(\beta)$  which is the iteration of gradient descent ( $\eta > 0$  is called the learning parameter). Show that at the minimizer of  $R_n$ ,  $\hat{\beta}_n$  we have that  $F_n(\hat{\beta}_n) = \hat{\beta}_n$ .

**Solution:**

In order to show that at the minimizer of  $R_n$ ,  $\hat{\beta}_n$  we have that  $F_n(\hat{\beta}_n) = \hat{\beta}_n$  we have to show that the gradient of the risk function is zero at the minimizer of  $R_n$ .

$$\frac{\partial R_n(\beta)}{\partial \beta} = \frac{\partial}{\partial \beta} \left( \eta \sum_{i=1}^n J(\beta) \right) = \eta \sum_{i=1}^n \frac{\partial}{\partial \beta} J(\beta) = \eta \sum_{i=1}^n x_i(e^{\beta^T x_i} - y_i)$$

To minimize this, we need to set the above equation to 0 in order to get the estimate  $\hat{\beta}_n$  i.e.

$$\eta \sum_{i=1}^n x_i(e^{\beta^T x_i} - y_i) = 0$$

$$e^{\beta^T x_i} - y_i = 0$$

Taking the log on both the sides we get

$$\hat{\beta}^T x_i - \log y_i = 0$$

Therefore, we have

$$\hat{\beta}^T = \log y_i x_i^{-1}$$

We know as a fact that since the  $R_n(\beta)$  is convex, the gradient of the risk function is zero at the minimizer of  $R_n$ . Therefore, we can conclude that the estimate  $\hat{\beta}_n$  is the minimizer of  $R_n$ .

At the minimizer of  $R_n$ , we have that  $F_n(\hat{\beta}_n) = \hat{\beta}_n$  as  $\nabla J_n = 0$  at  $\beta = \hat{\beta}_n$  as we can see from the above calculations.

1. I have a script to simulate from this model below. Implement the gradient descent algorithm above and show that with enough data ( $n$  large enough) the estimated  $\hat{\beta}$  approaches the true  $\beta$  (you can look at the sum of square error between these two vectors).

**Solution:**

```
In [1]: # Importing the required libraries for question 1
import numpy as np
import matplotlib.pyplot as plt

In [2]: ## Simulate from the Poisson regression model (use y,X)
np.random.seed(2022)
n, p = 1000, 20
X = np.random.normal(0,1,size = (n,p))
beta = np.random.normal(0,.,size = (p,1))
lamb = np.exp(X @ beta)
y = np.random.poisson(lamb)

In [3]: # Defining gradient descent function
def gradient_descent(X,y,eta,maxiter = 10000):
    """
    The gradient descent function takes in the X,y,eta and
    maxiter as inputs and returns the optimal beta.
    It initially starts from a random beta and tries to approximate
    F(beta) = beta - eta * grad J(beta)
    """
    Args:
        X (np.ndarray): n x p matrix of covariates
        y (np.ndarray): n x 1 vector of response
        eta (float): learning rate
        maxiter (int, optional): iterations for which the approximation is done.
        Defaults to 10000.
    """
    # Defining the initial beta
    random_beta = np.random.normal(0,0.2,size = (p))
    # Initializing the beta values
    beta = np.zeros((maxiter,20))
    # Getting the predictions
    predictions = np.exp(X @ random_beta)
    # Computing the gradient descent beta for the ith iteration
    random_beta = random_beta - eta * 1/(X.shape[0]) * X.T @ (predictions - y)
    # Storing the beta values
    beta[i+1,:] = random_beta
    # Storing the all obtained betas
    return beta

In [4]: # Checking for various learning rates that might be suitable for the
# problem
eta = [10, 0.1, 0.01, 0.001, 0.0001, 1e-05]
# Getting the mse for all etas
overall_mse = []
for eta in eta:
    # Getting the estimates
    estimated_beta = gradient_descent(X,y,eta)
    # Storing the mses at every iteration
    mse = []
    for estimate in estimated_beta:
        mse.append((estimate - beta)**2).mean()
    # Appending the mse for the particular eta
    overall_mse.append(mse)

In [5]: # Getting a plot for every eta value and seeing how quickly does
# the mse decrease
plt.figure(figsize=(10,5))
plt.plot(overall_mse)
plt.xlabel('Iterations')
plt.ylabel('MSE')
plt.legend()
plt.savefig('mse.png')

Out[5]: <matplotlib.legend.Legend at 0x1e7a3ee760>
```

From this above, we can conclude that the gradient descent algorithm converges to the true  $\beta$  as the learning rate  $\eta$  is small enough. In this case, the  $\eta$  seems to be  $\eta = 0.0001$  where it converges very smoothly between 0 and 2000 iterations.

## Exercise 2 (Regression and OLS) (35 pts, 5 pts each)

Consider the regression setting in which  $x_i \in \mathbb{R}^p$  and  $y_i \in \mathbb{R}$  for  $i = 1, \dots, n$ .

1. For a given regressor, let  $\hat{y}_i$  be prediction given  $x_i$ , and  $\hat{y}$  be the vector form. Show that linear regression can be written in the form

$$\hat{y} = H\hat{\beta}$$

where  $H$  is dependent on  $X$  (the matrix of where each row is  $x_i$ ), assuming that  $p < n$  and  $X$  is full rank. Give an expression for  $H$  or an expression for computing  $H$ .

**Solution:**

Given  $\hat{y} = H\hat{\beta}$ , in general for any  $p < n$ , we can write the linear regression as:

The predictions  $\hat{y}$  are given by  $X\hat{\beta}$  where  $X$  is the matrix of where each row is  $x_i$  and  $\hat{\beta}$  is the vector of estimated coefficients.

$$\hat{y} = X\hat{\beta}$$

By the OLS solution we know that  $\hat{\beta} = (X^T X)^{-1} X^T y$ . Substituting the above expression for  $\hat{\beta}$  we get:

$$\hat{y} = X\hat{\beta} = X(X^T X)^{-1} X^T y$$

To find the  $\hat{\beta}$ , we can use the OLS solution, where we try to minimize the following:

$$\begin{aligned} \arg \min_{\hat{\beta}} \|y - X\hat{\beta}\|_2^2 \\ = \arg \min_{\hat{\beta}} \|y - X\hat{\beta}\|_2^2 \\ = \arg \min_{\hat{\beta}} \|y - 2\hat{\beta}^T X^T y + \hat{\beta}^T X^T X \hat{\beta}\|_2^2 \end{aligned}$$

On taking the derivative of the above expression, we get:

$$\hat{\beta} = (X^T X)^{-1} X^T y$$

Correlating with the given expression  $\hat{y} = H\hat{\beta}$ , we get:

$$H = X(X^T X)^{-1} X^T$$

1. Assuming  $p < n$  and  $X$  is full rank, let  $U = UDV^T$  be the thin singular value decomposition where  $U$  is  $n \times p$ , and  $V$ ,  $D$  is  $p \times p$  ( $D$  is diagonal).

- a) Derive an expression for the OLS coefficients  $\hat{\beta} = \beta$  such that  $A$  is  $p \times p$  and depends on  $V$  and  $D$ , and  $b$  is a  $p$  vector and does not depend on  $D$ .

**Solution:** We know that the OLS coefficients are given by:

$$\hat{\beta} = (X^T X)^{-1} X^T y$$

Multiplying by  $U^T$  on both the sides and by property of SVD we know that  $U^T U = I$ , we get

$$\begin{aligned} U^T (UDV^T) y &= (U^T U) D V^T y = D V^T y = U^T y \\ \text{Since } D \text{ is a singular matrix, we can take the inverse of } D \text{ and multiply it by } I \text{ to get the OLS coefficients.} \\ V^T \hat{\beta} &= D^{-1} U^T y \\ V^T \hat{\beta} &= V D^{-1} U^T y \\ \hat{\beta} &= V D^{-1} U^T y \end{aligned}$$

On comparing the obtained expression with the given expression  $\hat{\beta} = A\hat{\beta}$ , we can say that,

$$A = V V^T = V D^{-1} U^T U = V D^{-1} U^T$$

- b) Describe a efficient method that precomputes these quantities separately

**Solution:**

SVD function can be used to compute the SVD of a matrix. In order to compute the above  $A$  and  $B$  we can write the following

```
In [6]: # Defining the library for question 2
from numpy.linalg import svd

In [7]: # Defining a function to calculate the A and B values
def custom_svd(X,y):
    # Getting the U,D,V values from the X matrix
    U,D,V = svd(X, full_matrices = False)
    # Getting inverse of D
    D_inv = np.diag(1/D)
    # Getting the A and B values
    A = U.T @ V @ D_inv
    b = U.T @ y
    return A,b
```

- c) Use the simulated data  $y$  and  $X$  in below to find  $\hat{\beta}$  using SVD.

```
In [8]: ## Simulate from the linear regression model (use y,X)
np.random.seed(2022)
n, p = 100, 20
X = np.random.normal(0,1,size = (n,p))
beta = np.random.normal(0,.,size = (p,1))
sigma = 1
y = np.random.normal(X @ beta, sigma**2)

In [9]: # Calculating the estimated beta using the above function
A,b = custom_svd(X,y)
beta_hat = A @ b
beta_hat
```

```
Out[9]: array([-0.17292927,  0.2403934 , -0.44231575, -0.03407534,  0.23607844,
        -0.12489189,  0.50945368,  0.03948573,  0.21520295, -0.03151865,
        -0.02374895,  0.13474846, -0.37652418, -0.01026382, -0.33508302,
        -0.04318508,  0.23211108, -0.20268112, -0.05872932, -0.24552763])
```

- d) Call a new data  $\tilde{X} \in \mathbb{R}^{m \times p}$ , derive an expression for the predicted  $\tilde{y}$  with  $\tilde{X}$  using SVD.

**Solution:**

We know that the predicted  $\tilde{y}$  is given by:

$$\tilde{y} = \tilde{X} \hat{\beta} \text{ where } \hat{\beta} = (X^T X)^{-1} X^T y$$

On substituting the SVD decomposition of  $\hat{\beta}$  we get:

$$\begin{aligned} \tilde{y} &= \tilde{X} \hat{\beta} = \tilde{X} V D^{-1} U^T y \\ &= U D V^T \tilde{X} V D^{-1} U^T y \end{aligned}$$

On simplifying the above expression we get,

$$\tilde{y} = U U^T y$$

1. Consider a regressor that performs OLS using the SVD above, but every instance of  $D$  will only use the largest  $r$  values on the diagonal, the remainder will be set to 0. Call this new  $\tilde{p} \times p$  matrix  $D_r$  ( $r < p$ ). Then the new coefficient vector is the OLS computed as if the design matrix is modified by  $X \rightarrow U D_r V^T$ .

- a) Given that you have computed  $b$  already, how could you make a method `change_rank` that recomputes  $A$  with  $D_r$  instead of  $D$ ?

**Solution:** Using the same implementation as before, we can modify the singular value decomposition to only use the largest  $r$  values on the diagonal. The function `change_rank` will be used to recompute the OLS coefficients with a lower rank using the same computation as `estimate_using_svd` function written in the previous part.

```
In [10]: # Defining the change_rank function
def change_rank(X,y,r = np.linalg.matrix_rank(X)):
    # Getting SVD decomposition
    U,D,V = svd(X, full_matrices = False)
    # svd() returns the singular values as an array so we convert it to a diagonal matrix
    D_inv = 1/D
    # getting r values and setting rest to 0
    D_inv = np.diag(D_inv[0:r], D_inv.shape[0] - (D_inv - r))
    D_inv = np.diag(D_inv)
    # Getting the A and B values
    A = U.T @ V @ D_inv
    b = U.T @ y
    # returning A and b
    return A,b
```

- b) Choose  $r = 10$ , recompute  $\hat{\beta}$  (call it  $\hat{\beta}_{lowRank}$ ) in Question 2-c.

**Solution:**

```
In [11]: # Defining A and B based on r = 10
A,b = change_rank(X,y,r = 10)
estimated_beta = A @ b
```

```
Out[11]: array([ 0.03112502, -0.04080854, -0.16394756,  0.03728632,  0.03755968,
        -0.04632204, -0.04193316,  0.13449484, -0.02375678, -0.02205015,
        -0.12489189,  0.02394537, -0.04746329, -0.01710239, -0.03823045,
        -0.02374895,  0.13474846, -0.37652418, -0.01026382, -0.33508302,
        -0.04318508,  0.23211108, -0.20268112, -0.05872932, -0.24552763])
```

## Exercise 3 (Subset selection) (15 pts)

Recall the subset selection problem with tuning parameter  $k$ .

$$\min_{\beta: |\beta|_0 \leq k} \|y - X\beta\|_2^2$$

where  $|\beta|_0 = |\{j = 1, \dots, p: \beta_j \neq 0\}|$ .

Notice that we can write this as

$$\min_{\beta: \text{supp}(\beta) \leq k} \|y - X\beta\|_2^2$$

where  $\text{supp}(\beta) = \{j = 1, \dots, p: \beta_j \neq 0\}$  ( $\text{supp}(\beta)$  is the support of  $\beta$ ).

1. (5 points) Write the subset selection problem in the following form

$$\min_{S \subseteq \{1, \dots, p\}, |S| \leq k} y^T P_S y,$$

where  $P_S$  is a projection.

**Solution:**

We know that the predictions are given by:

$$\hat{y} = X\hat{\beta} = H\hat{\beta} \text{ where } H = X(X^T X)^{-1} X^T$$

In order to minimize the above problem, we need to find a  $\hat{\beta}$  that would minimize the difference between  $y$  and  $\hat{y}$ , which happens to be  $\hat{\beta}$ . If we substitute  $\hat{\beta} = \hat{\beta}$ , then we get the new  $\hat{\beta}$  as mentioned above hence the problem further converts into:

$$\min_{S \subseteq \{1, \dots, p\}, |S| \leq k} \|y - H_S \hat{\beta}\|_2^2 = \min_{S \subseteq \{1, \dots, p\}, |S| \leq k} \|y - H_S y\|_2^2$$

Note that in the above expression, instead of the normal  $H$  which would consider all the predictors,  $H_S$  would be selecting all the predictors that are in the support of  $\beta$  which are in the set  $S$ .

On expanding the expression, we get:

$$\min_{S \subseteq \{1, \dots, p\}, |S| \leq k} y^T (I - H_S) y$$

Since  $H_S$  is a projection, we can also say that  $(I - H_S)$  is symmetric, hence we can write the above expression as:

$$\min_{S \subseteq \{1, \dots, p\}, |S| \leq k} y^T (I - H_S) y$$

We can substitute the expression for  $I - H_S$  as  $P_S$  and we get:

$$\min_{S \subseteq \{1, \dots, p\}, |S| \leq k} y^T P_S y$$

1. (10 points) Suppose that we have a nested sequence of models  $S_1 \subseteq S_2 \subseteq \dots \subseteq S_p$  such that  $|S_k| = k$  ( $|S_k|$  is the cardinality of  $S_k$  meaning that it contains  $k$  variables). Prove that

$$y^T P_{S_p} y \leq y^T P_{S_{k+1}} y$$

for  $k = 1, \dots, p-1$ . What does this tell us about the solution to the subset selection problem and the constraint  $|S| \leq k$ ?

(Hint: using the fact that  $X^T X$  is positive definite, write  $X^T X = V D V^T$ )

**Solution:**

Using the given inequality, we have:

$$\begin{aligned} y^T P_{S_k} y &\geq y^T P_{S_{k+1}} y \Rightarrow y^T (I - H_{S_k}) y \geq y^T (I - H_{S_{k+1}}) y \\ y^T y - y^T H_{S_k} y &\geq y^T y - y^T H_{S_{k+1}} y \\ y^T H_{S_k} y &\leq y^T H_{S_{k+1}} y \end{aligned}$$

Substituting the expression for  $H_S$  in its SVD form we get the following:

$$\begin{aligned} y^T U D V^T U D^{-1} U^T y &\leq y^T U D V^T U D^{-1} U^T y \\ y^T U D D^{-1} U^T y &\leq y^T U D D^{-1} U^T y \end{aligned}$$

Noticing that  $D D^{-1}$  is  $k \times k$  on the diagonal and rest 0's, we notice that the expression is in the form of  $\sum_{i=1}^k \lambda_i^2 x_i^2$  which can be simplified into:

$$\sum_{i=1}^k \lambda_i^2 x_i^2 \leq \sum_{i=1}^{k+1} \lambda_i^2 x_i^2$$

We can see the RHS is the summation of 1 extra term than on the LHS. Hence the inequality is satisfied.

## Exercise 4 (Ridge, lasso and adaptive lasso) (40 pts, 5 pts each)

For this exercise, it may be helpful to use the `sklearn.linear_model` module. I have also included a plotting tool for making the lasso path in ESL.

1. Load the training and test data using the script below. Fit OLS on the training dataset and compute the test error. Throughout you do not need to compute an intercept but you should normalize the  $X$  (divide by the column norms).

**Solution:**

```
In [12]: # Importing all the required packages for question 4
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error
from sklearn.linear_model import lars_path
from sklearn.preprocessing import normalize

In [13]: # Loading our data
import pickle
with open('hw1_data', 'rb') as f:
    y_tr, X_tr, y_te, X_te = pickle.load(f)

In [14]: # Normalizing the data by dividing the column by its norm
X_tr_scaled = normalize(X_tr, axis = 0)

In [15]: # Using the Linear Regression model from sci-kit learn
ols = LinearRegression(fit_intercept = False)
# Fitting the model
ols.fit(X_tr_scaled, y_tr)

Out[15]: LinearRegression(fit_intercept=False)
```

```
In [16]: # Getting the testing accuracy
predictions = ols.predict(X_te)
mse = mean_squared_error(y_te, predictions)
```

```
In [17]: # Getting the MSE
ols_sse = mean_squared_error(y_te, predictions)
print("The MSE for the OLS model is: {:.4f}".format(ols_sse))

The MSE for the OLS model is: 5209.3163
```

1. Ridge regression:

- a) Train and tune ridge regression using a validation set (choose LOOCV) and compute the test error (square error loss).

**Solution:**

```
In [17]: # Importing the required packages
from sklearn.linear_model import RidgeCV
from sklearn.model_selection import KFold

In [18]: # RidgeCV by default uses the LOO cross validation
# By giving it no alphas, it assumes default to be [0.1, 1, 10]
loocv_ridge = RidgeCV(store_cv_values = True, fit_intercept = False)
# Fitting the model
loocv_ridge.fit(X_tr_scaled, y_tr)

Out[18]: RidgeCV(alphas=array([ 0.1, 1., 10.]), fit_intercept=False,
```

```
In [19]: # Getting the predictions
loocv_ridge_predictions = loocv_ridge.predict(X_te)

# Getting the SSE for the RidgeCV model
loocv_ridge_sse = mean_squared_error(y_te, loocv_ridge_predictions)
print("The MSE for the RidgeCV model is: {:.4f}".format(loocv_ridge_sse))

The MSE for the RidgeCV model is: 3651.7683
```

- b) Repeat a) but using K-fold (you can choose  $K = 5$  or 10) cross validation, compute the test error. Compare the result to a). Comment on what you found.

**Solution:**

```
In [20]: # RidgeCV by default uses the LOO cross validation but by
# setting cv = 10, we get a 10-fold cross validation
# By giving it no alphas, it assumes default to be [0.1, 1, 10]
kfold_ridge = RidgeCV(store_cv_values = False, fit_intercept = False)
# Fitting the model
kfold_ridge.fit(X_tr_scaled, y_tr)

Out[20]: RidgeCV(alphas=array([ 0.1, 1., 10.]), cv=10, fit_intercept=False)
```

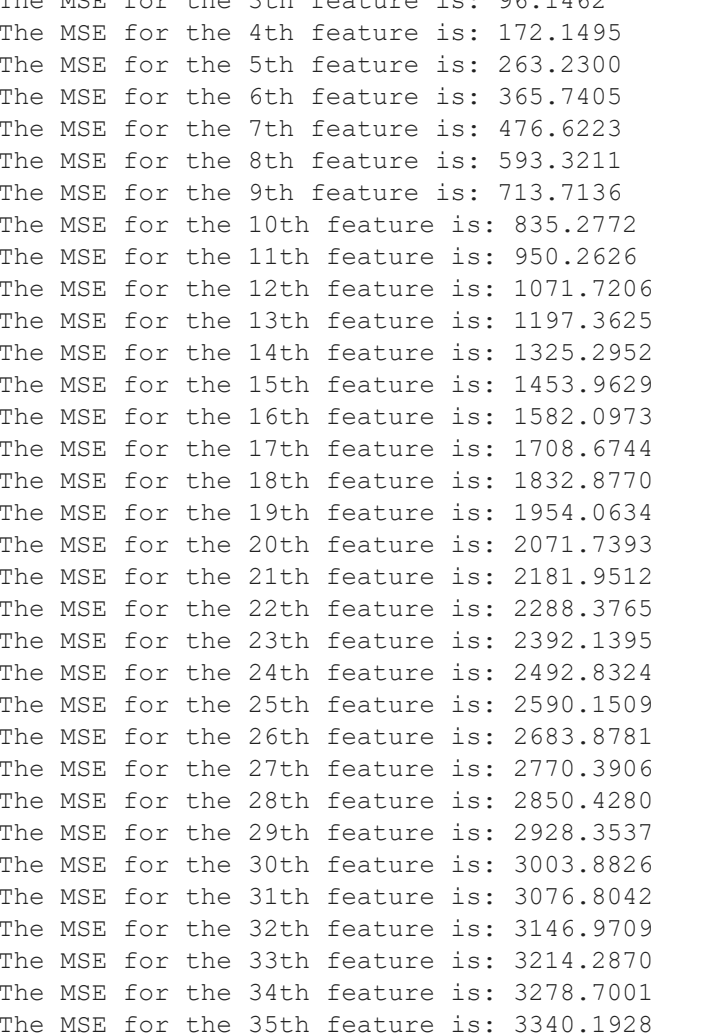
```
In [21]: # Getting the predictions
kfold_predictions = kfold_ridge.predict(X_te)


```



```
The MSE for the 0th feature is: 6.0565
The MSE for the 1th feature is: 628.4130
The MSE for the 2th feature is: 2103.2714
The MSE for the 3th feature is: 2733.1611
The MSE for the 4th feature is: 2641.3993
The MSE for the 5th feature is: 3675.6728
The MSE for the 6th feature is: 3835.6989
The MSE for the 7th feature is: 3661.7864
The MSE for the 8th feature is: 3961.7864
The MSE for the 9th feature is: 3970.7311
The MSE for the 10th feature is: 3972.9201
The MSE for the 11th feature is: 3973.3568
The MSE for the 12th feature is: 3975.9204
The MSE for the 13th feature is: 3980.9021
The MSE for the 14th feature is: 3991.6420
The MSE for the 15th feature is: 4006.7771
The MSE for the 16th feature is: 4016.6011
The MSE for the 17th feature is: 4028.9334
The MSE for the 18th feature is: 4072.1554
The MSE for the 19th feature is: 4094.7907
The MSE for the 20th feature is: 4107.4824
The MSE for the 21th feature is: 4110.1380
The MSE for the 22th feature is: 4121.5182
The MSE for the 23th feature is: 4117.3531
The MSE for the 24th feature is: 4135.6905
The MSE for the 25th feature is: 4135.6905
The MSE for the 26th feature is: 4138.9837
The MSE for the 27th feature is: 4166.8116
The MSE for the 28th feature is: 4173.5011
The MSE for the 29th feature is: 4178.9758
The MSE for the 30th feature is: 4181.9081
The MSE for the 31th feature is: 4193.5584
The MSE for the 32th feature is: 4205.5185
The MSE for the 33th feature is: 4210.4216
The MSE for the 34th feature is: 4237.7393
The MSE for the 35th feature is: 4262.0260
The MSE for the 36th feature is: 4264.7013
The MSE for the 37th feature is: 4285.8409
The MSE for the 38th feature is: 4306.4476
The MSE for the 39th feature is: 4300.8689
The MSE for the 40th feature is: 4313.1343
The MSE for the 41th feature is: 4359.4986
The MSE for the 42th feature is: 4316.6237
The MSE for the 43th feature is: 4322.9194
The MSE for the 44th feature is: 4348.0643
The MSE for the 45th feature is: 4350.0849
The MSE for the 46th feature is: 4362.9931
The MSE for the 47th feature is: 4392.5744
The MSE for the 48th feature is: 4399.8778
The MSE for the 49th feature is: 4402.1606
The MSE for the 50th feature is: 4405.8317
The MSE for the 51th feature is: 4417.8611
The MSE for the 52th feature is: 4421.5575
The MSE for the 53th feature is: 4439.5782
The MSE for the 54th feature is: 4439.8755
The MSE for the 55th feature is: 4442.2732
The MSE for the 56th feature is: 4455.9866
The MSE for the 57th feature is: 4456.9883
The MSE for the 58th feature is: 4476.3749
The MSE for the 59th feature is: 4477.0502
The MSE for the 60th feature is: 4497.0902
The MSE for the 61th feature is: 4498.0514
The MSE for the 62th feature is: 4512.2826
The MSE for the 63th feature is: 4512.8922
The MSE for the 64th feature is: 4529.1880
The MSE for the 65th feature is: 4539.0866
The MSE for the 66th feature is: 4557.6328
The MSE for the 67th feature is: 4581.3009
The MSE for the 68th feature is: 4591.6098
The MSE for the 69th feature is: 4598.7044
The MSE for the 70th feature is: 4636.2545
The MSE for the 71th feature is: 4655.8631
The MSE for the 72th feature is: 4665.9601
The MSE for the 73th feature is: 4668.3863
The MSE for the 74th feature is: 4774.6408
The MSE for the 75th feature is: 4781.8850
The MSE for the 76th feature is: 4794.5643
The MSE for the 77th feature is: 4806.4787
The MSE for the 78th feature is: 4819.9300
The MSE for the 79th feature is: 4825.3582
The MSE for the 80th feature is: 4831.2708
The MSE for the 81th feature is: 4848.5988
The MSE for the 82th feature is: 4892.3624
The MSE for the 83th feature is: 4892.2706
The MSE for the 84th feature is: 4892.7454
The MSE for the 85th feature is: 4895.9787
The MSE for the 86th feature is: 4987.5400
The MSE for the 87th feature is: 5003.4284
The MSE for the 88th feature is: 5060.2320
The MSE for the 89th feature is: 5088.9763
The MSE for the 90th feature is: 5094.8199
The MSE for the 91th feature is: 5144.4027
The MSE for the 92th feature is: 5200.7131
The MSE for the 93th feature is: 5212.8489
The MSE for the 94th feature is: 5221.6186
The MSE for the 95th feature is: 5226.5567
The MSE for the 96th feature is: 5251.4288
The MSE for the 97th feature is: 5288.5955
The MSE for the 98th feature is: 5299.0952
The MSE for the 99th feature is: 5354.7887
The MSE for the 100th feature is: 5355.6189
The MSE for the 101th feature is: 5366.6738
```

```
In [27]: plot_lars(coefs=coefs,lines=False, title = "LARS path")
```



On comparing the performance of the lasso path and the lars path, we see that the lasso path is more accurate than the lars path, having a lower MSE. However, the lars path converges quicker than the lasso path as seen where it takes 102 iterations to converge whereas the lars path takes only 101 iterations.

Next, LARS path seems to take variables into consideration even after they hit 0 due to soft thresholding, however in Lasso path, the variables are not considered. **Ex:** Red path in the both plots.

```
In [28]: # Seeing how many elements are common in the active sets
len(set(active_lasso) & set(active_lars))
```

Out[28]: 100

```
In [29]: # Doing an order by order comparison of the two paths
[1, 5, 8, 0, 23, 58, 33, 32]
```

Out[29]: [5, 8, 0, 23, 58, 33, 32]

From the above, we can see that although the elements are the same, the order of consideration into the model seems to be different. The above shown estimates are the only one those terms that match in terms of the order of the entry in the model.

1. Fit the lasso path with coordinate descent to the data. Compare the lasso path using coordinate descent with the path using lars. Comment on what you found.

**Solution:**

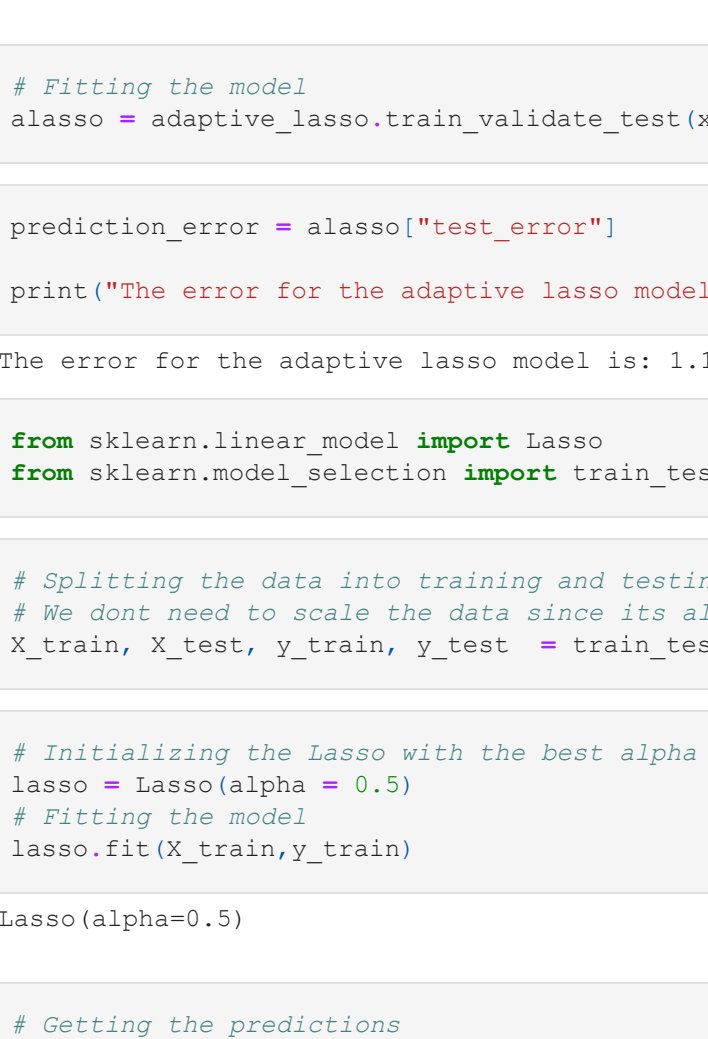
```
In [30]: # Using coordinate descent to get the coefficients
# Importing the required packages
from sklearn.linear_model import lasso_path
```

```
In [31]: # Getting the lasso path for training data using
# coordinate descent.
# By default, the lasso path is calculated by coordinate descent
alpha,coefs,_ = lasso_path(X_tr_scaled,y_tr)
```

```
# Calculating the MSE
for i in range(coefs.shape[1]):
    print("The MSE for the {}th feature is: {:.4f}".format(i,
        mean_squared_error(y_te,X_te @ coefs[:,i])))
```

```
The MSE for the 0th feature is: 6.0565
The MSE for the 1th feature is: 7.5029
The MSE for the 2th feature is: 39.5533
The MSE for the 3th feature is: 96.1462
The MSE for the 4th feature is: 172.1462
The MSE for the 5th feature is: 263.2300
The MSE for the 6th feature is: 365.7405
The MSE for the 7th feature is: 476.6223
The MSE for the 8th feature is: 593.3211
The MSE for the 9th feature is: 713.7136
The MSE for the 10th feature is: 835.2712
The MSE for the 11th feature is: 950.2626
The MSE for the 12th feature is: 1071.7206
The MSE for the 13th feature is: 1197.2625
The MSE for the 14th feature is: 1325.2952
The MSE for the 15th feature is: 1453.9629
The MSE for the 16th feature is: 1582.0973
The MSE for the 17th feature is: 1708.6744
The MSE for the 18th feature is: 1832.8770
The MSE for the 19th feature is: 1954.1634
The MSE for the 20th feature is: 2071.7393
The MSE for the 21th feature is: 2181.9512
The MSE for the 22th feature is: 2298.3765
The MSE for the 23th feature is: 2392.1395
The MSE for the 24th feature is: 2492.8324
The MSE for the 25th feature is: 2590.1509
The MSE for the 26th feature is: 2683.8781
The MSE for the 27th feature is: 2770.3906
The MSE for the 28th feature is: 2850.4280
The MSE for the 29th feature is: 2928.3537
The MSE for the 30th feature is: 3003.8826
The MSE for the 31th feature is: 3076.8142
The MSE for the 32th feature is: 3146.9709
The MSE for the 33th feature is: 3214.2870
The MSE for the 34th feature is: 3278.7001
The MSE for the 35th feature is: 3340.1928
The MSE for the 36th feature is: 3398.7771
The MSE for the 37th feature is: 3454.4878
The MSE for the 38th feature is: 3507.3788
The MSE for the 39th feature is: 3557.5182
The MSE for the 40th feature is: 3604.9856
The MSE for the 41th feature is: 3650.1262
The MSE for the 42th feature is: 3692.3519
The MSE for the 43th feature is: 3730.6918
The MSE for the 44th feature is: 3766.9864
The MSE for the 45th feature is: 3801.3030
The MSE for the 46th feature is: 3833.7501
The MSE for the 47th feature is: 3864.6366
The MSE for the 48th feature is: 3895.1212
The MSE for the 49th feature is: 3923.1824
The MSE for the 50th feature is: 3949.6969
The MSE for the 51th feature is: 3974.6577
The MSE for the 52th feature is: 3999.4135
The MSE for the 53th feature is: 4023.1580
The MSE for the 54th feature is: 4046.7250
The MSE for the 55th feature is: 4069.8404
The MSE for the 56th feature is: 4108.4212
The MSE for the 57th feature is: 4179.9616
The MSE for the 58th feature is: 4222.6148
The MSE for the 59th feature is: 4264.8383
The MSE for the 60th feature is: 4305.9015
The MSE for the 61th feature is: 4346.0474
The MSE for the 62th feature is: 4423.6967
The MSE for the 63th feature is: 4460.5803
The MSE for the 64th feature is: 4522.0205
The MSE for the 65th feature is: 4567.1863
The MSE for the 66th feature is: 4600.7201
The MSE for the 67th feature is: 4631.9071
The MSE for the 68th feature is: 4657.1260
The MSE for the 69th feature is: 4681.6366
The MSE for the 70th feature is: 4705.0444
The MSE for the 71th feature is: 4727.5015
The MSE for the 72th feature is: 4749.1967
The MSE for the 73th feature is: 4770.2507
The MSE for the 74th feature is: 4790.3261
The MSE for the 75th feature is: 4809.3367
The MSE for the 76th feature is: 4827.0700
The MSE for the 77th feature is: 4843.7920
The MSE for the 78th feature is: 4860.4992
The MSE for the 79th feature is: 4886.2425
The MSE for the 80th feature is: 4905.8637
The MSE for the 81th feature is: 4924.3930
The MSE for the 82th feature is: 4941.8927
The MSE for the 83th feature is: 4958.3353
The MSE for the 84th feature is: 4973.8215
The MSE for the 85th feature is: 4988.3099
The MSE for the 86th feature is: 5002.1131
The MSE for the 87th feature is: 5015.6558
The MSE for the 88th feature is: 5027.1632
The MSE for the 89th feature is: 5038.6658
The MSE for the 90th feature is: 5049.5298
The MSE for the 91th feature is: 5059.7348
The MSE for the 92th feature is: 5069.3021
The MSE for the 93th feature is: 5078.2704
The MSE for the 94th feature is: 5086.7121
The MSE for the 95th feature is: 5094.5884
The MSE for the 96th feature is: 5101.9585
The MSE for the 97th feature is: 5108.8930
```

```
In [32]: plot_lars(coefs=coefs,lines=False, title = "Lasso Path with coordinate descent")
```



On comparing the test error for Lasso using the LARS path and Lasso path, we get a lower test error for the Lasso path of 5108.89 as opposed to 5209.31

```
In [33]: coefs.shape
```

Out[33]: (100, 100)

1. Extract each active set from the lasso path and recompute the restricted OLS for each. Compute and compare the test error for each model.

```
In [34]: # For every model in the lasso path we look at
# all the non-zero coefficients
overall_mse = []
# X_te_scaled = normalize(X_te,axis = 0)
for coef in coefs.T:
    # Getting all non zero coefficients
    subset_idx = np.where(coef != 0)[0]
    # Getting a subset of the data
    X_tr_subset = X_tr_scaled[:,subset_idx]
    X_te_subset = X_te[:,subset_idx]
    if X_tr_subset.shape[1] == 0:
        continue
    # Fitting an OLS model
    restricted_ols = LinearRegression(fit_intercept = False)
    # Fitting the model
    restricted_ols.fit(X_tr_subset,y_tr)
    # Getting the predictions
    predictions = restricted_ols.predict(X_te_subset)
    # Getting the mse and storing it
    overall_mse.append(mean_squared_error(y_te,predictions))
```

```
In [35]: plt.plot(overall_mse)
```

Out[35]: [matplotlib.lines.Line2D at 0x1e7ee334a60]



We can see that introducing more and more variables is leading to an increase in the MSE for the model. This could be because the testing dataset is not well fitted.

1. Fit the lasso path with coordinate descent to the data. Compare the test error between the adaptive lasso and lasso for each returned coefficient. Comment on what you found.

```
In [36]: # Combining the training and testing data
X = np.concatenate((X_tr,X_te),axis = 0)
X_train,X_test = X[:n_train],X[n_train:]
y = np.concatenate((y_tr,y_te),axis = 0)
```

```
In [37]: # Importing the required libraries
import asgl
```

```
lambda1 = np.linspace(0,1,1000)
adaptive_lasso = asgl.TVJ(model = 'lm', penalization='alasso',lambda1 = lambda1,
    error_type='MSE',weight_technique='lasso', parallel = True,
    random_state=1,
    train_pct=0.5, validate_pct= .25)
```

```
In [38]: adaptive_lasso.alpha
```

Out[38]: 0.5

```
In [39]: # Fitting the model
lasso = adaptive_lasso.train_validate_test(x = X,y = y)
```

```
In [40]: prediction_error = lasso("test_error")
```

```
print("The error for the adaptive lasso model is: {:.4f}".format(prediction_error))
```

The error for the adaptive lasso model is: 1.1149

```
In [44]: from sklearn.linear_model import lasso
from sklearn.model_selection import train_test_split
```

```
In [45]: # Splitting the data into training and testing data
# We dont need to scale the data since its already scaled before
X_train,X_test,y_train,y_test = train_test_split(X, y, test_size=0.25, random_state=1)
```

```
In [46]: # Initializing the lasso with the best alpha in the adaptive lasso case
lasso = lasso(alpha = 0.5)
# Fitting the model
lasso.fit(X_train,y_train)
```

Out[46]: Lasso(alpha=0.5)

```
In [47]: # Getting the predictions
predictions = lasso.predict(X_test)
```

```
In [48]: # Getting the testing error
mean_squared_error(y_test,predictions)
```

Out[48]: 6.111467142644866

We can see that the error obtained in the adaptive lasso is lower than the lasso. We get an error of 1.07 for the adaptive lasso and 6.11 for the lasso.