

Advanced Machine Learning

CPSC 8420

Problem 1: Ridge Regression

$$\text{To prove: } (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I}_p)^{-1} \mathbf{A}^T = \mathbf{P}^T (\mathbf{A} \mathbf{A}^T + \lambda \mathbf{I}_n)^{-1}$$

det, L.H.S. = K.

R.H.S = P

$$(\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I}_p)^{-1} \mathbf{A}^T = \mathbf{K}$$

multiplying by $(\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I}_p)$

$$(\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I}_p) (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I}_p)^{-1} \mathbf{A}^T = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I}_p) \mathbf{K}$$

Multiplying both sides by $(\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I}_n)$

$$\mathbf{A}^T (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I}_n) = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I}_p) (\mathbf{K}) (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I}_n)$$

$$\mathbf{A}^T \mathbf{A} \mathbf{A}^T + \lambda \mathbf{A}^T \mathbf{I}_n = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I}_p) (\mathbf{K}) (\mathbf{A} \mathbf{A}^T + \lambda \mathbf{I}_n) \quad (1)$$

Now, for R.H.S:

$$\mathbf{A}^T (\mathbf{A} \mathbf{A}^T + \lambda \mathbf{I}_n)^{-1} = \mathbf{P}$$

similarly, we can say that:

$$(\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I}_p) \mathbf{A}^T (\mathbf{A} \mathbf{A}^T + \lambda \mathbf{I}_n)^{-1} (\mathbf{A} \mathbf{A}^T + \lambda \mathbf{I}_n) =$$

$$(\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I}_p) \mathbf{P} (\mathbf{A} \mathbf{A}^T + \lambda \mathbf{I}_n)$$

$$\Rightarrow (\mathbf{A}^T \mathbf{A} \mathbf{A}^T + \lambda \mathbf{I}_p) \mathbf{A}^T = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I}_p) \mathbf{P} (\mathbf{A} \mathbf{A}^T + \lambda \mathbf{I}_n)$$

$$A^T A A^T + \lambda A^T I_n = (A^T A + \lambda I_p)(P)(A A^T + \lambda I_n) \quad \textcircled{1}$$

Subtracting $\textcircled{2}$ from $\textcircled{1}$, we get:

$$(A^T A + \lambda I_p)(P)(A A^T + \lambda I_n) - (A^T A + \lambda I_p)(P)(A A^T + \lambda I_n) = (A^T A + \lambda I_p)(P)(A A^T + \lambda I_n) - (A^T A + \lambda I_p)(P)(A A^T + \lambda I_n)$$

Multiplying by $(A^T A + \lambda I_p)^{-1}$ on both sides, we get:

$$K(A A^T + \lambda I_p) = P(A A^T + \lambda I_n)$$

Similarly, multiplying by $(A A^T + \lambda I_n)$

we get $K = P \cdot A \cdot (A^T A + \lambda I_n)$

$$\text{LHS} = RHS$$

```
% Define the size of the random matrix (e.g., a 4x4 matrix)
n = 100;
p = [10, 100, 1000, 2000];
l = 0.1;
I_n = eye(n);
t_p = [];
t_n = [];

% Generate a random matrix of integers between 1 and 10
for i = 1:length(p)
    I_p = eye(p(i));
    A = randi([1, 100], n, p(i));
    tic;
    B = inv(A'*A + l .* I_p) * A';
    t_p(i) = toc;
    tic;
    C = A'*inv(A*A' + l * I_n);
    t_n(i) = toc;
end

t_p = 1000 * t_p;
t_n = 1000 * t_n;
disp("t_p: ")
```

t_p:

```
disp(t_p);
```

0.8866 0.4024 34.2900 158.3368

```
disp("t_n: ")
```

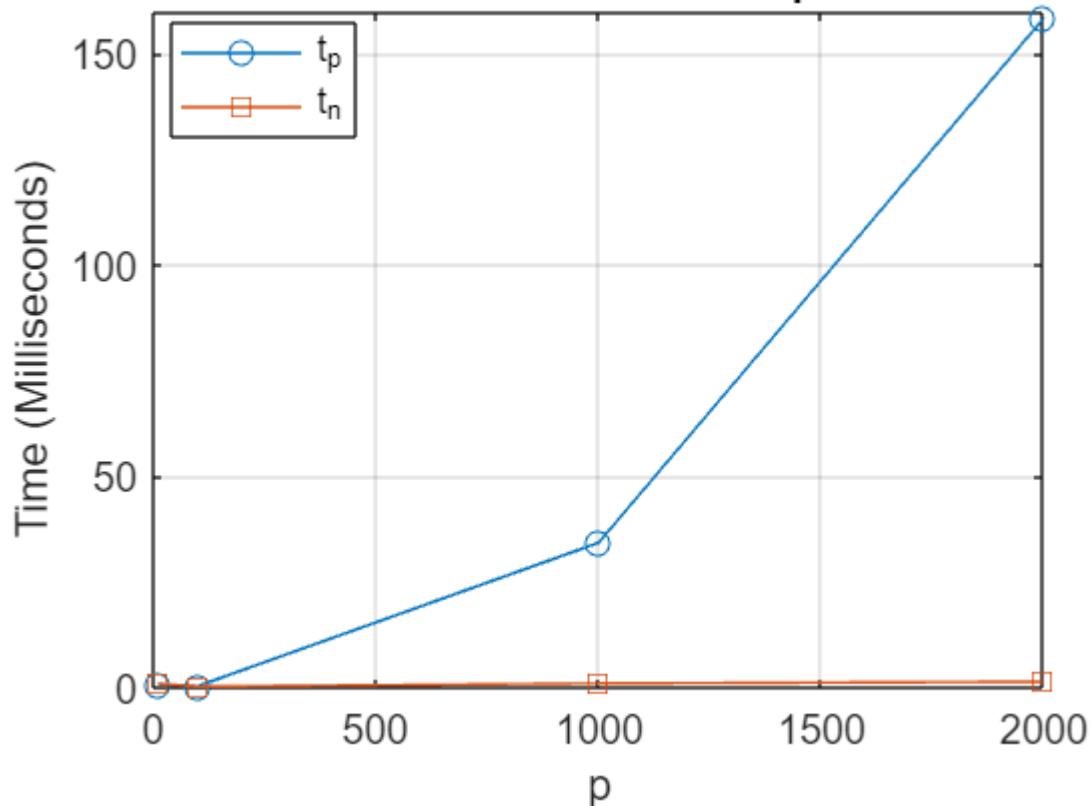
t_n:

```
disp(t_n);
```

0.9951 0.3015 1.0300 1.4470

```
% Create a plot
figure;
plot(p, t_p, '-o', 'DisplayName', 't_p');
hold on;
plot(p, t_n, '-s', 'DisplayName', 't_n');
xlabel('p');
ylabel('Time (Milliseconds)');
title('Execution Time vs. p');
legend('Location', 'Best');
grid on;
hold off;
```

Execution Time vs. p



We can see that the execution time t_n is considerably less than t_p . So, this form should be preferred for solving a regressor with large number of samples.

Ans 2) To solve: ~~min ||AX + XC + Y||_F^2~~ ~~order~~

$$(V - (X)) \min_x ||AX + XC + Y||_F^2 \quad \text{--- } ①$$

given: $\text{vec}(AXC) = (C^T \otimes A)\text{vec}(X) - \text{--- } ②$

Using: $\min ||\text{vec}(AX) + \text{vec}(XC) + \text{vec}(Y)||_F^2$

using eqn ②:

$$\text{vec}(AX) = \text{vec}(AXI) \quad \text{--- } ③$$

$$= (I^T \otimes A)\text{vec}(X)$$

$$\text{vec}(XC) = \text{vec}(IXC)$$

$$= (C^T \otimes I)\text{vec}(X) \quad \text{--- } ④$$

from ③; ④ we have :

$$\min_x ||(I^T \otimes A)\text{vec}(X) + (C^T \otimes I)\text{vec}(X) - Y||_F^2$$

$$\min_x ||[(I^T \otimes A) + (C^T \otimes I)]\text{vec}(X) - Y||_F^2$$

As, $(I^T \otimes A) + (C^T \otimes I) = M$

we have: $\min ||M \cdot \text{vec}(X) - Y||_F^2$

using least squares Approach \hat{y} (LSM)

$$\nabla f(\text{vec}(x)) = 2M^T(M \cdot \text{vec}(x) - y)$$

Putting $\nabla f(\text{vec}(x)) = 0$

$$2M^T(M \cdot \text{vec}(x) - y) = 0$$

$$M^T M \text{vec}(x) = M^T y$$

$$\boxed{\text{vec}(x) = (M^T M)^{-1} M^T y}$$

```
% Define the size of our example matrices
n = 2;

% Generate a random nxn matrix for A
A = randn(n, n);
disp('A: ')
```

A:

```
disp(A);
```

```
-1.0891    0.5525
 0.0326    1.1006
```

```
% Generate a random nxn matrix for X_star (optimal X)
X_star = randn(n, n);
disp('X_Star: ')
```

X_Star:

```
disp(X_star);
```

```
1.5442   -1.4916
 0.0859   -0.7423
```

```
% Generate a random nxn matrix for C
C = randn(n, n);
disp('C: ')
```

C:

```
disp(C);
```

```
-1.0616   -0.6156
 2.3505    0.7481
```

```
% Compute the matrix Y based on the relation: Y = AX* + X*C
Y = A * X_star + X_star * C;
```

```
% Create an identity matrix of size nxn
I = eye(n);
```

```
% Construct matrix M using the Kronecker product
% This is used for reshaping the matrix multiplication into vector form
M = kron(I, A) + kron(C.', I);
```

```
% Use matrix division (or the inverse of M) to compute vecX from Y
vecX = M \ Y(:);
```

```
% Reshape vecX back into matrix form to get X
X = reshape(vecX, [n, n]);
```

```
% Display the computed X and M matrices  
disp('X: ')
```

X:

```
disp(X);
```

```
1.5442 -1.4916  
0.0859 -0.7423
```

```
disp('M: ')
```

M:

```
disp(M);
```

```
-2.1506 0.5525 2.3505 0  
0.0326 0.0390 0 2.3505  
-0.6156 0 -0.3410 0.5525  
0 -0.6156 0.0326 1.8487
```

vec(X) and vec(X*) are equivalent.

Problem 3 $(\mathbf{X}^T \mathbf{A}) \mathbf{b}_{LS} = (\mathbf{X} \mathbf{A})_{LS}$

Vanilla Linear Regression:

$$(1) - f(\beta) = \min_{\beta} \|y - A\beta\|_2^2$$

using least squares approach to minimize $f(\beta)$:

$$\|\mathbf{Y} - (\mathbf{X}^T \mathbf{A} \beta)^T + \nabla f(\beta)\|_2$$

$$\nabla f(\beta)_2 = -2 \mathbf{A}^T (\mathbf{y} - \mathbf{A} \hat{\beta}_{LS})$$

$$-2 \mathbf{A}^T (\mathbf{y} - \mathbf{A} \hat{\beta}_{LS})_2 = 0$$

$$\mathbf{M} = (\mathbf{I} \otimes \mathbf{A}^T) + (\mathbf{A} \otimes \mathbf{I})$$

$$\mathbf{A}^T \cdot \mathbf{A} \hat{\beta}_{LS} = \mathbf{A}^T \cdot \mathbf{y}$$

$$\hat{\beta}_{LS} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}$$

$$\hat{\beta}_{LS} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}$$

Given $A^T A = I$, so $\boxed{\beta_{LS} = I A^T y}$

Ridge linear regression:

$$f(\beta) = \min_{\beta} \|y - A\beta\|_2^2 + \lambda \|\beta\|_2^2$$

using least squares approach to minimize $f(\beta)$

$$\nabla f(\beta) = 0$$

$$\nabla f(\beta) = -2A^T(y - A\beta) + 2\lambda\beta = 0$$

$$A^T(A\beta - y) + \lambda\beta = 0$$

$$A^T A \beta - A^T y + \lambda\beta = 0$$

$$(A^T A + \lambda I)\beta = A^T y$$

Given $A^T A = I$

$$\hat{\beta}_{\text{ridge}} = ((I + \lambda I)^{-1} A^T y)$$

$$\hat{\beta}_{\text{ridge}} = \frac{1}{1+\lambda} \cdot I \cdot A^T y$$

Lasso linear regression

$$f(\beta) = \min_{\beta} \|y - A\beta\|_2^2 + \lambda \|\beta\|_1$$

Given $A^T A = I$ (i.e. $A \rightarrow$ orthogonal)

for lasso problem, $\hat{\beta}_{\text{lasso}}$ satisfies the subgradient optimality condition,

$$\nabla f(\beta) + \lambda \text{sgn}(\beta_{\text{lasso}}) = 0 \quad \text{--- (1)}$$

$$f_L(\beta) = \min_{\beta} \|A\beta - y\|_2^2$$

Using optimality conditions for lasso condⁿ

$$\nabla f_L(\beta) = A^T (A\beta - y) \\ \Rightarrow A^T A \cdot \beta = A^T y.$$

$$\text{given } A^T A \rightarrow I_A$$

$$\nabla f_L(\beta) = \beta - A^T y$$

$$\text{from (1) } \beta - A^T y + \lambda \text{sgn}(\beta) = 0$$

$$\hat{\beta}_{\text{lasso}} = A^T y - \lambda \text{sgn}(\hat{\beta}_{\text{lasso}})$$

Now the soft thresholding operator $S_\lambda(z)$ is defined as:

$$S_\lambda(z) = \text{sgn}(z) \max(|z| - \lambda, 0)$$

Given optimality condition each component $\hat{\beta}_{\text{Lasso}}^*(i)$ of vector $\hat{\beta}_{\text{Lasso}}$ can be found by applying the soft thresholding operator to vector $A^+ b$ by.

$$\hat{\beta}_{\text{Lasso}}^*(i) = s_2(A^T y)_i$$

Subset selection Method:

$$f(\beta) = \min_{\beta} \frac{1}{2} \|y - A\beta\|_2^2 + \lambda \|\beta\|_1.$$

Given $A^T A = I$, ($A \rightarrow$ orthogonal)
 This means that correlation b/w any two columns of A is zero.

∴ correlation b/w y & j^{th} column of A is just $A_j^T y$.

so, for each predictor j :
 compute $c_j = A_j^T y$

$$\beta_j^{\text{subset}} = \begin{cases} \beta_j^* > 0 \text{ if } |c_j| > \lambda \\ \beta_j^* = c_j \text{ otherwise} \end{cases}$$

```

data = load(['C:\Users\param\OneDrive - Clemson University' ...
    '\Desktop\CU\Fall23\Advance Machine Learning\HW1\housing.data']);

x = data(:, 1:13);
y = data(:,14);

% Vanilla Linear Regression
beta_ols = pinv(x' * x) * x' * y;

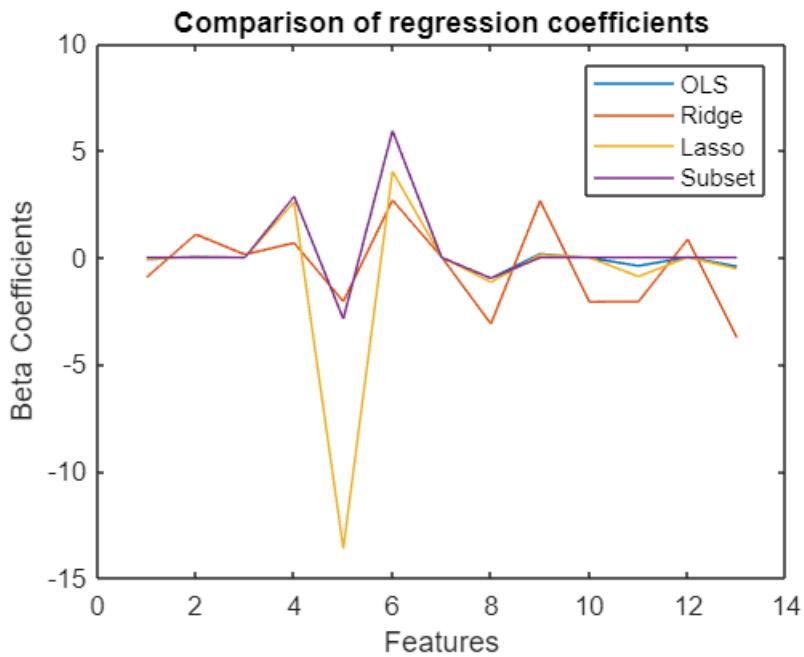
lambda = 0.1;
% Ridge regression
beta_ridge = ridge(y, x, lambda);

% Lasso regression
[beta_lasso, fitinfo] = lasso(x, y, 'Lambda', lambda, 'intercept', true);

% Subset (Simple implementation using a thresholding on OLS coefficients)
lambda_subset = 0.5;
beta_subset = beta_ols;
beta_subset(abs(beta_ols) < lambda_subset) = 0;

% Plot
figure;
plot(1:13, beta_ols, 'DisplayName', 'OLS');
hold on;
plot(1:13, beta_ridge, 'DisplayName', 'Ridge');
plot(1:13, beta_lasso, 'DisplayName', 'Lasso');
plot(1:13, beta_subset, 'DisplayName', 'Subset');
title('Comparison of regression coefficients');
xlabel('Features');
ylabel('Beta Coefficients');
legend();
hold off;

```



Problem 4: Linear Regression and it's regression

Preparing train and test sets for models

```
data = load('housing.data');
x = data(:, 1:13);
y = data(:, 14);
[n, d] = size(x);

seed = 2; rand('state', seed); randn('state', seed);
perm = randperm(n); % remove any possible ordering fx
x = x(perm, :); y = y(perm);
Ntrain = 300;
Xtrain = x(1:Ntrain, :); ytrain = y(1:Ntrain);
Xtest = x(Ntrain + 1:end, :); ytest = y(Ntrain + 1:end);
```

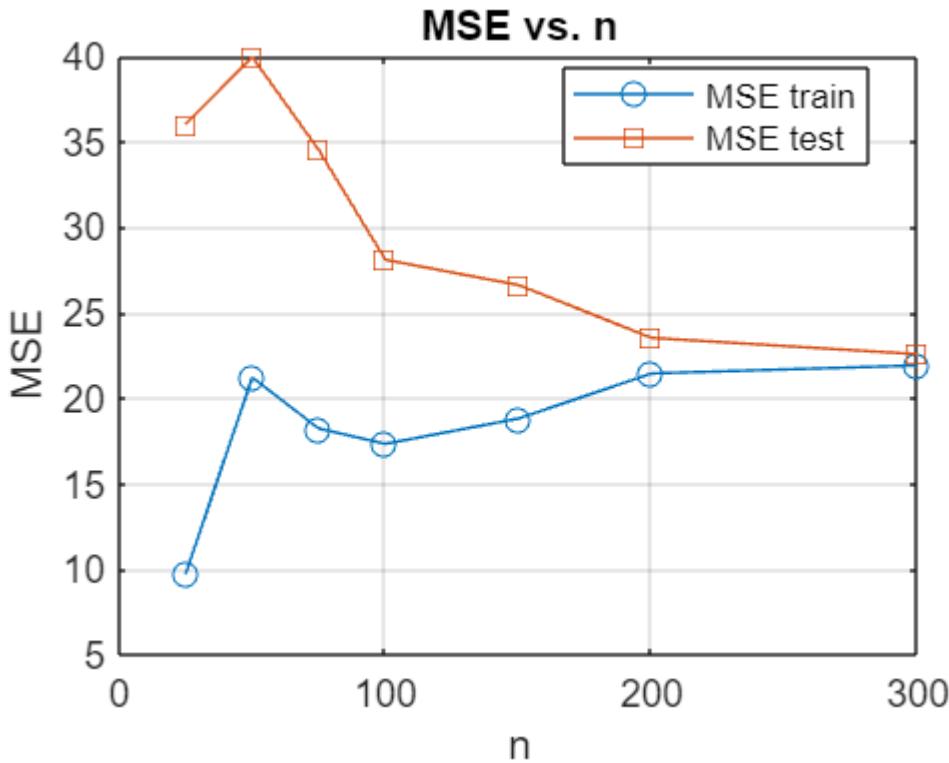
Part 2

```
n = [25 50 75 100 150 200 300];

MSE_train = zeros(size(n));
MSE_test = zeros(size(n));

for i = 1:length(n)
    % Normalizing the train data
    [st_array1, m1 , s1] = zscore(Xtrain(1:n(i), :));
    X_train = [ones(size(Xtrain(1:n(i), :), 1), 1) st_array1];
    % Normalizing the test data
    s_test = (Xtest-m1)./s1;
    X_test = [ones(size(Xtest, 1), 1) s_test];
    y_train = ytrain(1:n(i));
    weights = pinv(X_train'*X_train)*(X_train' * y_train);
    MSE_train(i) = ((X_train * weights - y_train)' * (X_train * weights - y_train))/n(i);
    MSE_test(i) = ((X_test * weights - ytest)'*(X_test * weights - ytest))/size(X_test, 1);
end

% Create a plot
figure;
plot(n, MSE_train, '-o', 'DisplayName', 'MSE train');
hold on;
plot(n, MSE_test, '-s', 'DisplayName', 'MSE test');
xlabel('n');
ylabel('MSE');
title('MSE vs. n');
legend('Location', 'Best');
grid on;
hold off;
```



Expanding the size of the training dataset equips the model with a richer source of information, typically resulting in superior adaptability to unseen data, consequently reducing test error.

In scenarios where the training dataset is limited in size, the model can excessively tailor its predictions to the confined data, occasionally veering into overfitting. With the enlargement of the training dataset, the information to be absorbed becomes more multifaceted and intricate to accommodate flawlessly, leading to an upturn in training error.

Furthermore, as the training dataset's dimensions augment and edge nearer to the complete dataset, the model undergoes exposure to a more comprehensive portion of the available data during training. As a result, its performance on both the training set and the test set starts to draw closer, culminating in the convergence of the two error curves.

Part 3

```

deg = [1 2 3 4 5 6];

MSE_train = zeros(size(deg));
MSE_test = zeros(size(deg));

for i = 1:length(deg)
    % Train Feature expansion
    X_train = degexpand(Xtrain, deg(i), 1);
    [X_train(:, 2:end), mean, std] = zscore(X_train(:, 2:end));
    y_train = ytrain;
    y_test = ytest;
    % Test feature expansion

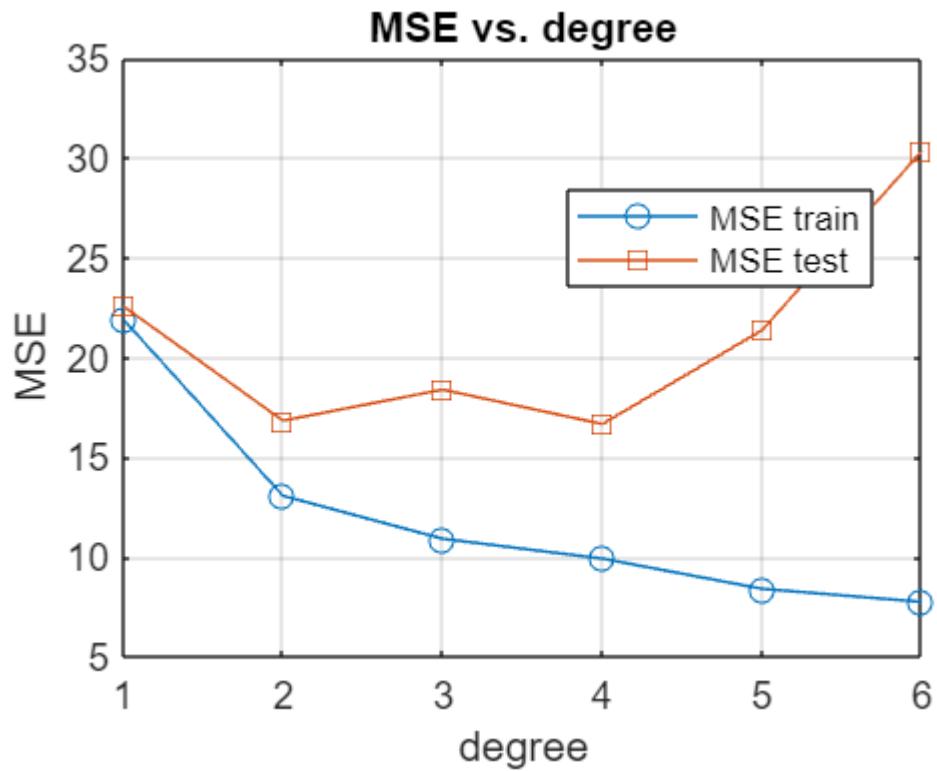
```

```

X_test = degexpand(Xtest, deg(i), 1);
X_test(:, 2:end) = (X_test(:, 2:end)-mean)./std;
weights = pinv(X_train'*X_train) * X_train' * y_train;
MSE_train(i) = ((X_train * weights - y_train)' * (X_train * weights - y_train))/size(X_train, 1);
MSE_test(i) = ((X_test * weights - y_test)'*(X_test * weights - y_test))/size(X_test, 1);
end

% Create a plot
figure;
plot(deg, MSE_train, '-o', 'DisplayName', 'MSE train');
hold on;
plot(deg, MSE_test, '-s', 'DisplayName', 'MSE test');
xlabel('degree');
ylabel('MSE');
title('MSE vs. degree');
legend('Location', 'Best');
grid on;
hold off;

```

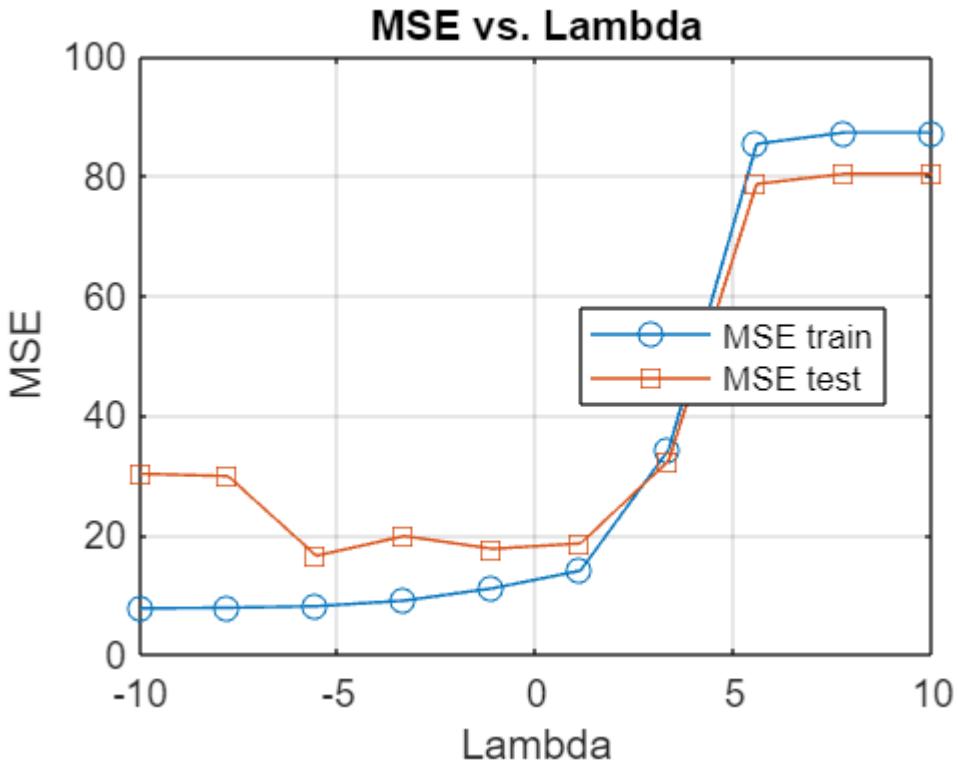


As the polynomial degree is raised, the model gains greater adaptability, allowing it to conform more closely to the training data, subsequently decreasing the training error.

In the early stages, the increased adaptability of the model aids in capturing the fundamental data patterns, which, in turn, results in a decrease in the test error. However, there comes a point where the model begins to overfit. In essence, it starts fitting the training data with excessive precision, encompassing not only the genuine patterns but also the noise and outliers. This diminishes its capacity to generalize effectively to new data, ultimately leading to an uptick in test error.

Part 4

```
lambdas = [0 logspace(-10, 10, 10)];  
  
X_train = degexpand(Xtrain, 6, 1);  
[X_train(:, 2:end), mean, std] = zscore(X_train(:, 2:end));  
y_train = ytrain;  
y_test = ytest;  
X_test = degexpand(Xtest, 6, 1);  
X_test(:, 2:end) = (X_test(:, 2:end)-mean)./std;  
I = eye(size(X_train, 2));  
% To avoid regularizing the bias term, put I(0,0) = 0  
I = [zeros(size(I, 1), 1), I(:,2:end)];  
  
MSE_train = zeros(size(lambdas));  
MSE_test = zeros(size(lambdas));  
  
for i = 1:length(lambdas)  
    weights = pinv(X_train'*X_train + lambdas(i)* I) * X_train' * y_train;  
    MSE_train(i) = (((X_train * weights - y_train)' * (X_train * weights - y_train)))/size(X_tr  
    MSE_test(i) = (((((X_test * weights - ytest)'*(X_test * weights - ytest))))/size(X_test, 1);  
end  
  
% Create a plot  
figure;  
plot(log10(lambdas), MSE_train, '-o', 'DisplayName', 'MSE train');  
hold on;  
plot(log10(lambdas), MSE_test, '-s', 'DisplayName', 'MSE test');  
xlabel('Lambda');  
ylabel('MSE');  
title('MSE vs. Lambda');  
legend('Location', 'Best');  
grid on;  
hold off;
```



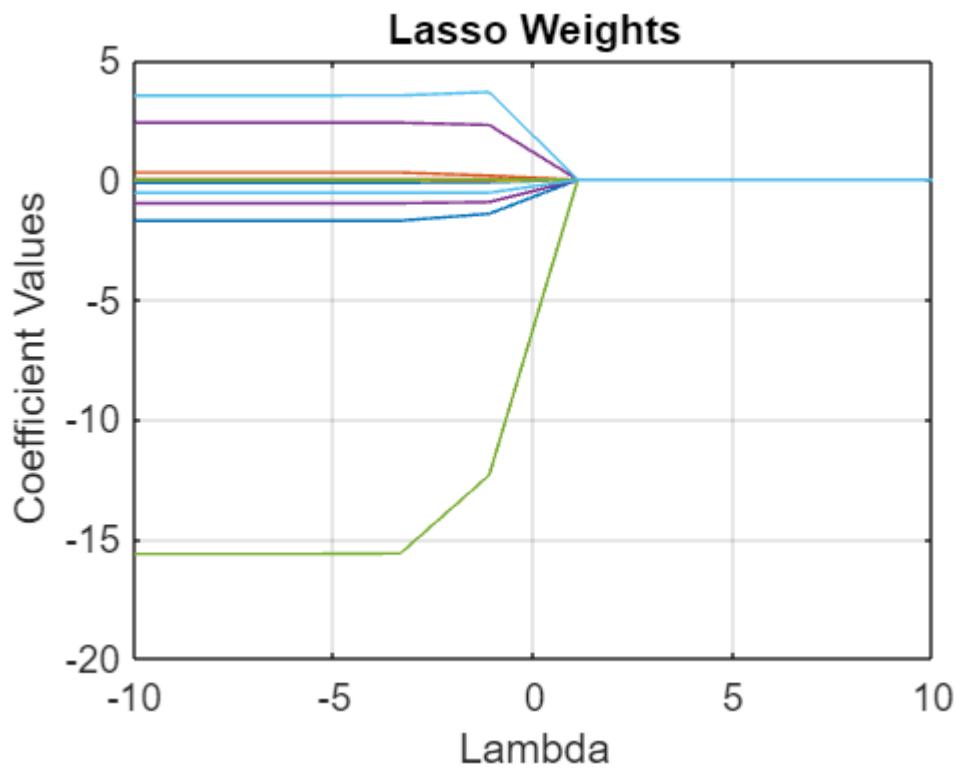
Initially, regularization helps in avoiding overfitting. But as seen from the plot, as the regularization increases, the model starts to underfit which leads to an increased mean squared error.

Part 5

```
lambdas = logspace(-10, 10, 10);

[weights, FitInfo] = lasso(Xtrain, ytrain, 'Lambda', lambdas);

figure;
plot(log10(lambdas), weights);
xlabel('Lambda');
ylabel('Coefficient Values');
title('Lasso Weights');
grid on;
```



We observe that as the regularization increases, the weights are pushed towards zero.