

COMP9417 - Machine Learning Homework 1: Regularized Regression & Numerical Optimization

a)

The equation outlining the explicit gradient update:

$$x^{(k+1)} = x^k - \alpha(A^T(Ax^k - b) + \gamma x^k)$$

A print out of the first 5 ($k = 5$ inclusive) and last 5 rows of iterations:

```
(cv) > comp9417 git:(master) python -u "/Users/yueyifei/Desktop/COMP9417/comp9417/hw1/code/q_a.py"
k = 0, x(k) = [1,1,1,1]
k = 1, x(k) = [0.98,0.98,0.98,0.98]
k = 2, x(k) = [0.9624,0.9804,0.9744,0.9584]
k = 3, x(k) = [0.9427,0.9824,0.9668,0.9433]
k = 4, x(k) = [0.9234,0.9866,0.9598,0.9295]
k = 5, x(k) = [0.9044,0.9916,0.9526,0.9169]
k = 272, x(k) = [0.0666,1.3366,0.4928,0.3251]
k = 273, x(k) = [0.0666,1.3366,0.4928,0.325]
k = 274, x(k) = [0.0665,1.3366,0.4927,0.325]
k = 275, x(k) = [0.0664,1.3367,0.4927,0.3249]
k = 276, x(k) = [0.0663,1.3367,0.4927,0.3249]
```

A screen shot of code for a):

```

comp9417 > hw1 > code > q_a.py > [o]
You, 1 hour ago | 1 author (You)
1 import numpy as np
2
3 alpha = 0.1
4 gamma = 0.2
5
6 A = np.array([[1,2,1,-1], [-1,1,0,2], [0,-1,-2,1]])      # A: 3 x 4
7 b = np.array([3, 2, -2])                                     # b: 3 x 1
8 A_t = A.T                                                 # A_t: 4 x 3
9 x0 = np.array([1,1,1,1]).reshape(4, 1)                      # x_0: 1 x 4 -> 4 x 1
10
11 x_list = []
12 x_list.append(x0)
13
14 x = x0
15 k = 1
16
17 while True:
18     grad = A_t @ (A @ x - b) + gamma * x
19     if np.linalg.norm(grad, ord=2) < 0.001:
20         break
21     new_x = x - alpha * grad
22     x = new_x
23     x_list.append(new_x)
24     k = k + 1
25
26 length = len(x_list)
27
28 for i in range(0, 6):
29     x = x_list[i].reshape(1, 4)[0]
30     print("k = {}, x(k) = [{} , {} , {} , {}]".format(i, round(x[0],4), round(x[1],4), round(x[2],4), round(x[3],4)))
31
32 for j in range(length - 5, length):
33     x = x_list[j].reshape(1, 4)[0]
34     print("k = {}, x(k) = [{} , {} , {} , {}]".format(j, round(x[0],4), round(x[1],4), round(x[2],4), round(x[3],4)))
35

```

b)

This termination condition means that when the gradient of the function is less than 0.001, we consider that the algorithm has converged to a minimizer of f.

If other conditions remain unchanged and the gradient required by the termination condition is limited to a smaller value, we will calculate a more accurate value of x that minimizes f. Therefore, we need to complete more iterations, and the time for the algorithm to converge to a minimizer of f also becomes longer.

c)

A print out of the first 5 (k = 5 inclusive) and last 5 rows of iterations:

```

(cv) > comp9417 git:(master) python -u "/Users/yueyifei/Desktop/COMP9417/comp9417/hw1/code/q_c.py"
k = 0, x(k) = [1.0,1.0,1.0,1.0]
k = 1, x(k) = [0.98,0.98,0.98,0.98]
k = 2, x(k) = [0.9624,0.9804,0.9744,0.9584]
k = 3, x(k) = [0.9427,0.9824,0.9668,0.9433]
k = 4, x(k) = [0.9234,0.9866,0.9598,0.9295]
k = 5, x(k) = [0.9044,0.9916,0.9526,0.9169]
k = 272, x(k) = [0.0666,1.3366,0.4928,0.3251]
k = 273, x(k) = [0.0666,1.3366,0.4928,0.325]
k = 274, x(k) = [0.0665,1.3366,0.4927,0.325]
k = 275, x(k) = [0.0664,1.3367,0.4927,0.3249]
k = 276, x(k) = [0.0663,1.3367,0.4927,0.3249]

```

A screen shot of code for c):

```

comp9417 > hw1 > code > q_c.py > [o] i
You, 4 hours ago | 1 author (You)
1 import torch
2 import torch.nn as nn
3 from torch import optim
4
5 gamma = 0.2
6 alpha = 0.1
7
8 A = torch.tensor([[1., 2., 1., -1.], [-1., 1., 0., 2.], [0., -1., -2., 1.]])
9 b = torch.tensor([[3.], [2.], [-2.]])
10 A_t = A.t()
11 x0 = torch.tensor([[1.], [1.], [1.], [1.]])
12
13 x_dict = {}
14 x_dict[0] = x0.tolist()
15 x = x0.clone()
16

You, 4 hours ago | 1 author (You)
17 class MyModel(nn.Module):
18     def __init__(self):
19         super().__init__()
20         self.x = nn.Parameter(x, requires_grad=True)
21     def forward(self, x):
22         z = torch.mul(torch.pow(torch.norm(torch.sub(torch.mm(A, x), b)), 2), 1/2)
23         w = torch.mul(torch.pow(torch.norm(x), 2), gamma/2)
24         return torch.add(z, w)
25
26 model = MyModel()
27 optimizer = optim.SGD(model.parameters(), lr=alpha)
28 terminationCond = False
29
30 k = 1
31
32 while not terminationCond:
33     optimizer.zero_grad()
34     loss = model(model.x)
35     loss.backward()
36     optimizer.step()
37
38 if torch.norm(model.x.grad) < 0.001:
39     terminationCond = True
40 else:
41     x_dict[k] = model.x.tolist()
42     k = k + 1
43
44 length = len(x_dict)
45
46 for i in range(0, 6):
47     x = x_dict[i]
48     print("k = {}, x(k) = {{}, {}, {}, {}}}".format(i, round(x[0][0], 4), round(x[1][0], 4), round(x[2][0], 4), round(x[3][0], 4)))
49
50 for j in range(length - 5, length):
51     x = x_dict[j]
52     print("k = {}, x(k) = {{}, {}, {}, {}}}".format(j, round(x[0][0], 4), round(x[1][0], 4), round(x[2][0], 4), round(x[3][0], 4)))

```

d)

We apply fit_transform() on the training data to scale the training data (standardize the remaining features) and also learn the scaling parameters of that data which are the mean and variance of the features of the training set. After centering the target variable, we can create our training set and testing set and output the result.

A print out of the means and variances of features and the first and last rows of the 4 requested objects:

```
(base) → hw1 git:(master) ✘ python -u "/Users/yueyifei/Desktop/COMP9417/comp9417/hw1/code/q_degh.py"
The means of features:
[ 3.81916720e-16  3.55271368e-17  2.66453526e-17  1.59872116e-16
 -6.21724894e-17  1.28785871e-16 -1.33226763e-16]

The variances of features:
[1. 1. 1. 1. 1. 1.]

The first and last rows of X_train:
[ 0.85045499  0.15536099  0.65717702  0.07581929  0.17782345 -0.69978222
  1.18444912]
[-0.1942498   0.6920138  -0.24615909  0.47665602  0.43155489  0.65991828
  0.03820804]

The first and last rows of X_test:
[ 1.24221929  0.83512122 -0.99893918  0.57176982  1.27732635  0.53630914
 -0.72595268]
[ 0.58927879 -1.13260576 -0.99893918 -1.61584759  0.17782345 -0.26715025
  0.80236876]

The first and last rows of Y_train:
[2.003675]
[-1.076325]

The first and last rows of Y_test:
[-1.936325]
[2.213675]
```

A screen shot of code for d):

```

comp9417 > hw1 > code > q_degh.py > ...
You, 56 seconds ago | 1 author (You)
 1  from sklearn.preprocessing import StandardScaler
 2  from sklearn.model_selection import train_test_split
 3  import pandas as pd
 4  import numpy as np
 5  import matplotlib.pyplot as plt
 6  import copy
 7
 8
 9  ### Code for Question d
10  data = pd.read_csv('CarSeats.csv')
11
12  label = ['Sales']
13  target = data[label]
14
15  categorical_features = ['ShelveLoc', 'Urban', 'US']
16
17  numerical_data = data.drop(label + categorical_features, axis=1)
18
19  scaler = StandardScaler()
20  numerical_data_scaled = scaler.fit_transform(numerical_data)
21
22  print('The means of features:')
23  print(numerical_data_scaled.mean(axis=0))
24  print()
25  print('The variances of features:')
26  print(numerical_data_scaled.var(axis=0))
27  print()
28
29  target_centered = (target - target.mean(axis=0)).to_numpy()
30  X_train, X_test, Y_train, Y_test = train_test_split(numerical_data_scaled, target_centered, test_size=0.5, shuffle=False)
31
32  print('The first and last rows of X_train:')
33  print(X_train[0])
34  print(X_train[-1])
35  print()
36  print('The first and last rows of X_test:')
37  print(X_test[0])
38  print(X_test[-1])
39  print()
40  print('The first and last rows of Y_train:')
41  print(Y_train[0])
42  print(Y_train[-1])
43  print()
44  print('The first and last rows of Y_test:')
45  print(Y_test[0])
46  print(Y_test[-1])
47  print()

```

e)

Since $\hat{\beta}_{\text{Ridge}} = \underset{\beta}{\operatorname{argmin}} \frac{1}{n} \|y - X\beta\|_2^2 + \phi \|\beta\|_2^2$,

then $\hat{\beta}_{\text{Ridge}} = \underset{\beta}{\operatorname{argmin}} \frac{1}{n} (y - X\beta)^T (y - X\beta) + \phi \|\beta\|^2$.

Therefore, to get closed form expression, we have

$$\begin{aligned} & \frac{1}{n} (2(y - X\beta)(-X)) + 2\phi\beta = 0 \\ \Rightarrow & (y - X\beta)(-X) + n\phi\beta = 0 \\ \Rightarrow & -yX + X^T X\beta + n\phi\beta = 0 \\ \Rightarrow & (X^T X + \cancel{n\phi})\beta = X^T y \\ \therefore & \hat{\beta}_{\text{Ridge}} = (X^T X + n\phi I)^{-1} X^T y \end{aligned}$$

A print out of the value of the ridge solution:

```
The value of the ridge solution based on X_train and Y_train:
[[ 0.680673 ]
 [ 0.28229334]
 [ 0.65157017]
 [ 0.00834835]
 [-1.17129533]
 [-0.400892 ]
 [-0.10063355]]
```

A screen shot of code for e):

```
50  ### Code for Question e
51  phi = 0.5
52  size = len(data.columns) - len(label + categorical_features)
53  X_t = X_train.T
54
55  beta_hat = np.linalg.inv(X_t @ X_train + phi * len(X_train) * np.identity(size)) @ X_t @ Y_train
56
57  print('The value of the ridge solution based on X_train and Y_train:')
58  print(beta_hat)
```

f)

$$\begin{aligned} \text{Since } L(\beta) &= \frac{1}{n} \|y - X\beta\|_2^2 + \phi \|\beta\|_2^2 \\ &= \frac{1}{n} \sum_{i=1}^n (y_i - x_i \beta)^2 + \phi \sum_{j=1}^m \beta_j^2, \end{aligned}$$

~~therefore~~ and $L(\beta) = \frac{1}{n} \sum_{i=1}^n L_i(\beta)$.

~~Therefore~~ Then we have,

$$\begin{aligned} n L(\beta) &= \sum_{i=1}^n (y_i - x_i \beta)^2 + n \phi \sum_{j=1}^m \beta_j^2 \\ &= \sum_{i=1}^n L_i(\beta), \end{aligned}$$

then we have

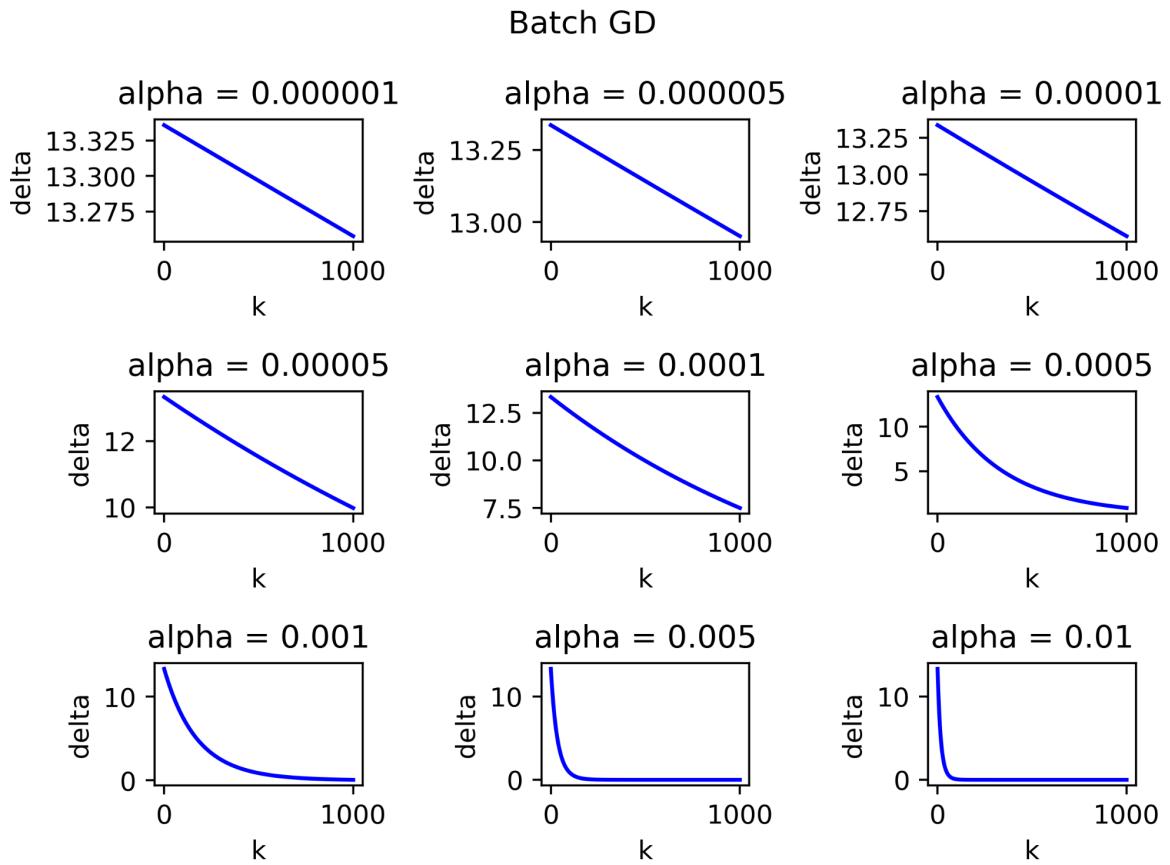
$$L_i(\beta) = (y_i - x_i \beta)^2 + \phi \sum_{j \neq i} \beta_j^2$$

Therefore, $L_n(\beta) = (y_n - x_n \beta)^2 + \phi \sum_{j \neq n} \beta_j^2$

$$\begin{aligned} \nabla L_n(\beta) &= \frac{\partial L_n(\beta)}{\partial \beta} = 2(y_n - x_n \beta)(-x_n) + 2\phi \beta \\ &= -2(y_n - x_n \beta)x_n + 2\phi \beta \end{aligned}$$

Therefore, $\nabla L_n(\beta) = -2(y_n - x_n \beta)x_n + 2\phi \beta$

g)



We choose the step-size $\alpha = 0.005$, in fact, both 0.005 and 0.01 have good effects, which can make the algorithm converge faster. But we choose 0.005 to further avoid the potential error caused by large step size.

A print out of the train and test MSE requested:

The train MSE is : 4.559027077966185
The test MSE is : 4.380486336397683

A screen shot of code for g):

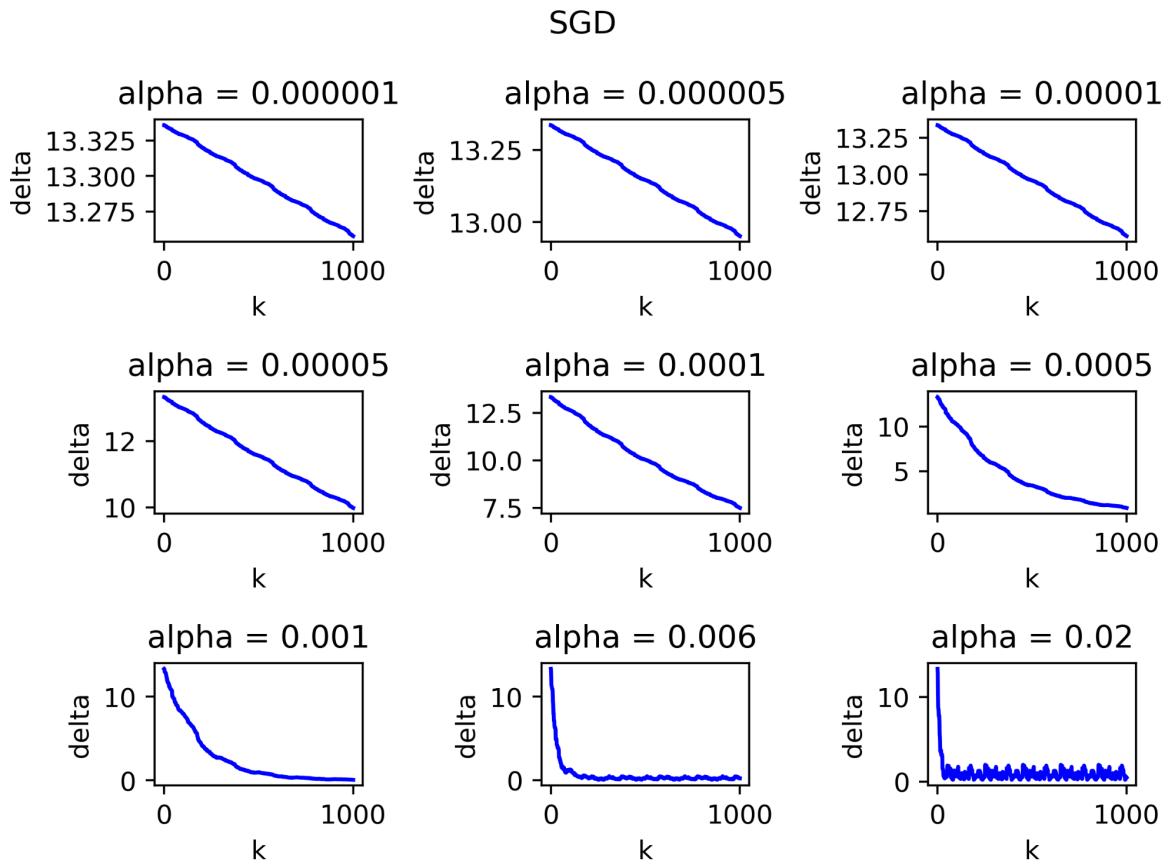
```

61  ### Code for Question g
62  beta_0 = np.ones(X_train.shape[1]).reshape(7, 1)
63
64  phi = 0.5
65  epochs = 1000
66  alphas = [0.000001, 0.000005, 0.00001, 0.00005, 0.0001, 0.0005, 0.001, 0.005, 0.01]
67
68  def loss(X_train, Y_train, beta, phi):
69      loss_beta = 1/len(X_train) * (np.linalg.norm(Y_train - X_train @ beta, ord=2) ** 2) + phi * (np.linalg.norm(beta, ord=2) ** 2)
70  return loss_beta
71
72  def grad(X_train, Y_train, beta, phi):
73      grad_beta = -2 * (1/len(X_train)) * X_train.T @ (Y_train - X_train @ beta) + 2 * phi * beta
74  return grad_beta
75
76  def get_gd_betas(beta_0, X_train, Y_train, alpha, phi, epochs):
77      beta = beta_0
78      betas = []
79      betas.append(beta)
80      for k in range(epochs):
81          b_updated = beta - alpha * grad(X_train, Y_train, beta, phi)
82          beta = b_updated
83          betas.append(b_updated)
84  return betas
85
86  def get_gd_deltas(beta_0, X_train, Y_train, alpha, phi, epochs):
87      deltas = []
88      betas = get_gd_betas(beta_0, X_train, Y_train, alpha, phi, epochs)
89      loss_hat = loss(X_train, Y_train, beta_hat, phi)
90      for beta in betas:
91          delta = loss(X_train, Y_train, beta, phi) - loss_hat
92          deltas.append(delta)
93  return deltas
94
95  i = 1
96  for alpha in alphas:
97      deltas = get_gd_deltas(beta_0, X_train, Y_train, alpha, phi, epochs)
98      axes = plt.subplot(3, 3, i)
99      plt.plot(range(epochs + 1), deltas, color='blue')
100     plt.xlabel('k')
101     plt.ylabel('delta')
102     plt.title('alpha = ' + str(np.format_float_positional(alpha)))
103     i = i + 1
104 plt.suptitle('Batch GD')
105 plt.tight_layout()
106 plt.savefig('GD_plot.png', dpi=400)
107 plt.show()

108  beta_best = get_gd_betas(beta_0, X_train, Y_train, 0.005, phi, epochs)[-1]
109
110  train_MSE = 1/len(X_train) * (np.linalg.norm(Y_train - X_train @ beta_best) ** 2)
111  print('The train MSE is :', train_MSE)
112
113  test_MSE = 1/len(X_train) * (np.linalg.norm(Y_test - X_test @ beta_best) ** 2)
114  print('The test MSE is :', test_MSE)
115
116

```

h)



We choose the step-size $\alpha = 0.006$, because when the step size used by the algorithm is smaller than this value, although the descent process is stable, but the convergence speed is slow, and when the step size is larger than this, the descending process is very unstable. Therefore, we choose $\alpha = 0.006$.

A print out of the train and test MSE requested:

**The train MSE is : 4.661749176041872
The test MSE is : 4.447228989660326**

A screen shot of code for h :

```

118  ### Code for Question h
119  beta_0 = np.ones(X_train.shape[1]).reshape(7, 1)
120  phi = 0.5
121  epochs = 5
122  sgd_alphas = [0.000001, 0.000005, 0.00001, 0.00005, 0.0001, 0.0005, 0.001, 0.006, 0.02]
123
124  def get_sgd_betas(beta_0, X_train, Y_train, alpha, phi, epochs):
125      beta = beta_0
126      betas = []
127      betas.append(beta)
128      for i in range(epochs):
129          for j in range(len(X_train)):
130              x = X_train[j]
131              y = Y_train[j]
132              x = x.reshape(7,1)
133              x_t = x.T
134              b_updated = beta - alpha * (-2 * x @ (y - x_t @ beta) + 2 * phi * beta)
135              beta = b_updated
136              betas.append(b_updated)
137      return betas
138
139  def get_sgd_deltas(beta_0, X_train, Y_train, alpha, phi, epochs):
140      deltas = []
141      betas = get_sgd_betas(beta_0, X_train, Y_train, alpha, phi, epochs)
142      loss_hat = loss(X_train, Y_train, beta_hat, phi)
143      for beta in betas:
144          delta = loss(X_train, Y_train, beta, phi) - loss_hat
145          deltas.append(delta)
146      return deltas
147
148  i = 1
149  for alpha in sgd_alphas:
150      deltas = get_sgd_deltas(beta_0, X_train, Y_train, alpha, phi, epochs)
151      axes = plt.subplot(3, 3, i)
152      plt.plot(range(epochs * len(X_train) + 1), deltas, color='blue')
153      plt.xlabel('k')
154      plt.ylabel('delta')
155      plt.title('alpha = ' + str(np.format_float_positional(alpha)))
156      i = i + 1
157  plt.suptitle('SGD')
158  plt.tight_layout()
159  plt.savefig('SGD_plot.png', dpi=400)
160  plt.show()
161
162  beta_best = get_sgd_betas(beta_0, X_train, Y_train, 0.006, phi, epochs)[-1]
163
164  train_MSE = 1/len(X_train) * (np.linalg.norm(Y_train - X_train @ beta_best) ** 2)
165  print('The train MSE is :', train_MSE)
166
167  test_MSE = 1/len(X_train) * (np.linalg.norm(Y_test - X_test @ beta_best) ** 2)
168  print('The test MSE is :', test_MSE)
169

```

i)

According to our results, the descending process of SGD has some fluctuation and is not very stable, while the descending process of GD is more stable. There is no big difference between the two in terms of the convergence speed, so we choose to use GD.

j)

$$\begin{aligned} \text{Since } L(\beta) &= \frac{1}{n} \|y - X\beta\|_2^2 + \phi \|\beta\|_2^2 \\ &= \frac{1}{n} \|y - X_j\beta_j - X_{-j}\beta_{-j}\|_2^2 + \phi \|\beta\|_2^2, \end{aligned}$$

then $\frac{\partial L(\beta)}{\partial \beta_j} = \frac{2}{n} (y - X_j\beta_j - X_{-j}\beta_{-j})(-X_j) + 2\phi \beta_j.$

Let $\frac{\partial L(\beta)}{\partial \beta_j} = 0$

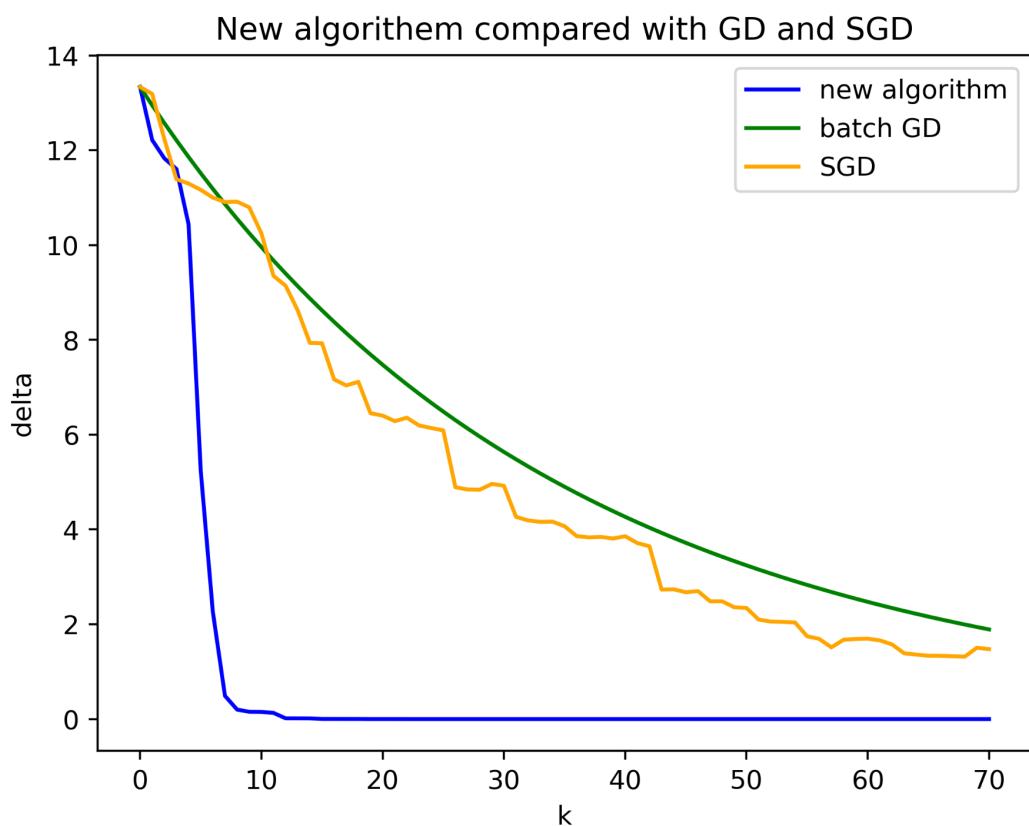
$$\begin{aligned} \frac{2}{n} (y - X_j\hat{\beta}_j - X_{-j}\beta_{-j})(-X_j) + 2\phi \hat{\beta}_j &= 0 \\ \hat{\beta}_j &= \frac{x_j y - x_j x_{-j} \beta_{-j}}{x_j x_j + n\phi} \end{aligned}$$

Therefore,

$$\hat{\beta}_1 = \frac{x_1 y - x_1 x_{-1} \beta_{-1}}{x_1 x_1 + n\phi},$$

$$\hat{\beta}_j = \frac{x_j y - x_j x_{-j} \beta_{-j}}{x_j x_j + n\phi} \quad \text{for } j = 1, 2, \dots, p.$$

k)



A print out of the train and test MSE requested:

The train MSE is : 4.5589067305440665
The test MSE is : 4.380429183660091

A screen shot of code for k):

```

171  ### Code for Question k
172  cycles = 10
173  p = len(X_train[0])
174  beta_0 = np.ones(X_train.shape[1]).reshape(7, 1)
175  phi = 0.5
176
177
178 def get_new_betas(beta_0, X_train, Y_train, phi):
179     beta = beta_0
180     betas = []
181     betas.append(beta)
182     y = Y_train
183     for i in range(cycles):
184         for j in range(p):
185             x_j = X_train[:, j].reshape(1, len(X_train))
186             x_nj = np.delete(X_train, j, 1)
187             beta_nj = np.delete(beta, j, 0)
188             bj_updated = (x_j @ y - x_j @ x_nj @ beta_nj) / (x_j @ x_j.T + len(X_train) * phi)
189             beta = copy.deepcopy(beta)
190             beta[j] = bj_updated
191             betas.append(beta)
192     return betas
193
194
195 def get_new_deltas(beta_0, X_train, Y_train, phi):
196     deltas = []
197     betas = get_new_betas(beta_0, X_train, Y_train, phi)
198     loss_hat = loss(X_train, Y_train, beta_hat, phi)
199     for beta in betas:
200         delta = loss(X_train, Y_train, beta, phi) - loss_hat
201         deltas.append(delta)
202     return deltas
203
204
205 deltas = get_new_deltas(beta_0, X_train, Y_train, phi)
206
207 plt.plot(range(cycles * p + 1), deltas, color='blue', label='new algorithm')
208 plt.plot(range(cycles * p + 1), get_gd_deltas(beta_0, X_train, Y_train, 0.005, phi, 1000)[0:cycles * p + 1], color='green', label='batch GD')
209 plt.plot(range(cycles * p + 1), get_sgd_deltas(beta_0, X_train, Y_train, 0.006, phi, 5)[0:cycles * p + 1], color='orange', label='SGD')
210 plt.xlabel('k')
211 plt.legend(loc='upper right') # creates legend in top right corner of plot
212 plt.ylabel('delta')
213 plt.title('New algorithem compared with GD and SGD')
214 plt.savefig('comparison.png', dpi=400)
215 plt.show()
216
217 beta_best = get_new_betas(beta_0, X_train, Y_train, phi)[-1]
218
219 train_MSE = 1/len(X_train) * (np.linalg.norm(Y_train - X_train @ beta_best) ** 2)
220 print('The train MSE is :', train_MSE)
221
222 test_MSE = 1/len(X_train) * (np.linalg.norm(Y_test - X_test @ beta_best) ** 2)
223 print('The test MSE is :', test_MSE)

```

l)

Less reliable. If we also standardize the target value, the proportion of the target value predicted by the model is different. The model needs to transform the predicted target value to the original value before the standardization, and there may be errors in the process.