

COMP9417 - Machine Learning Homework 2: Newton's Method and Mean Squared Error of Estimators

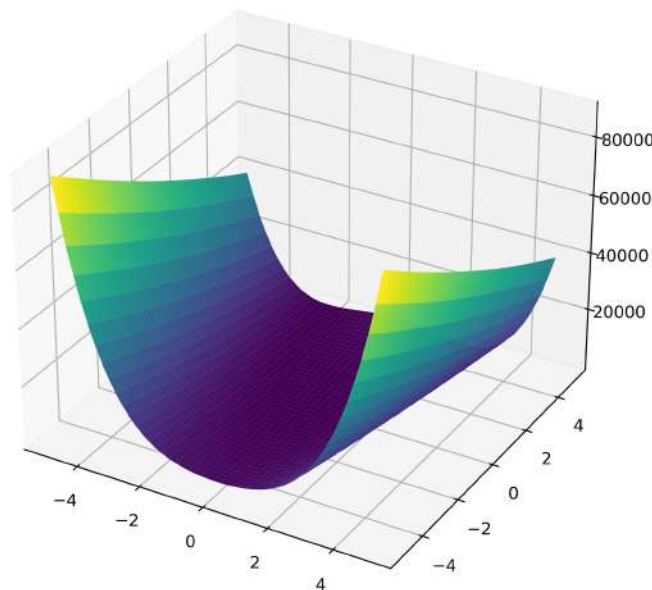
Question 1

a)

When the input x^k is a vector, then the first derivative of x^k is $\nabla f(x^k)$, and the second derivative of x^k is Hessian matrix which can be calculated as $H(x^k) = \nabla(\nabla f(x^k))$. According to the update rule equation 1, we have $x^{k+1} = x^k - \frac{\nabla f(x^k)}{\nabla(\nabla f(x^k))}$ by placing the first and second derivatives in the same position. Therefore, we have $x^{k+1} = x^k - (H(x^k))^{-1} \nabla f(x^k)$ which is a generalization of equation 1 to functions with vector inputs.

b)

The plot:



The code used to generate the plot:

```
q1_b.py > ...
1  import matplotlib.pyplot as plt
2  from mpl_toolkits import mplot3d
3  import numpy as np
4
5
6  def func(x, y):
7      return 100 * (y - x ** 2) ** 2 + (1 - x) ** 2
8
9
10 # create two one-dimensional grids using linspace
11 x = np.linspace(-5, 5, 50)
12 y = np.linspace(-5, 5, 50)
13 # combine the two one-dimensional grids into one two-dimensional grid
14 X, Y = np.meshgrid(x, y)
15
16 Z = func(X, Y)
17
18 fig = plt.figure(figsize=(7, 7))
19 ax = fig.add_subplot(projection='3d')
20 ax.plot_surface(X, Y, Z, cmap='viridis')
21 plt.savefig('qb.png', dpi=400)
22 plt.show()
23
```

The gradient and Hessian of f :

$$b) f(x, y) = 100(y - x^2)^2 + (1 - x)^2$$

$$\begin{aligned} dx = \frac{\partial f(x, y)}{\partial x} &= \frac{\partial}{\partial x} (100(y - x^2)^2) + \frac{\partial}{\partial x} ((1 - x)^2) \\ &= -400x(y - x^2) - 2(1 - x) \end{aligned}$$

$$\begin{aligned} dy = \frac{\partial f(x, y)}{\partial y} &= \frac{\partial}{\partial y} (100(y - x^2)^2) + \frac{\partial}{\partial y} ((1 - x)^2) \\ &= 200(y - x^2) \end{aligned}$$

$$\therefore \text{gradient} = \nabla f(x, y) = (-400x(y - x^2) - 2(1 - x), 200(y - x^2))^T$$

$$\begin{aligned} h_{00} = \frac{\partial^2 f(x, y)}{\partial x^2} &= \frac{\partial dx}{\partial x} = \frac{\partial (-400x(y - x^2))}{\partial x} + \frac{\partial (-2(1 - x))}{\partial x} \\ &= -400(y - 3x^2) + 2 \end{aligned}$$

$$\begin{aligned} h_{01} = \frac{\partial^2 f(x, y)}{\partial x \partial y} &= \frac{dx}{\partial y} = \frac{\partial (-400x(y - x^2))}{\partial y} + \frac{\partial (-2(1 - x))}{\partial y} \\ &= -400x \end{aligned}$$

$$\begin{aligned} h_{10} = \frac{\partial^2 f(x, y)}{\partial y \partial x} &= \frac{dy}{\partial x} = \frac{\partial (200(y - x^2))}{\partial x} \\ &= -400x \end{aligned}$$

$$\begin{aligned} h_{11} = \frac{\partial^2 f(x, y)}{\partial y^2} &= \frac{dy}{\partial y} = \frac{\partial (200(y - x^2))}{\partial y} \\ &= 200 \end{aligned}$$

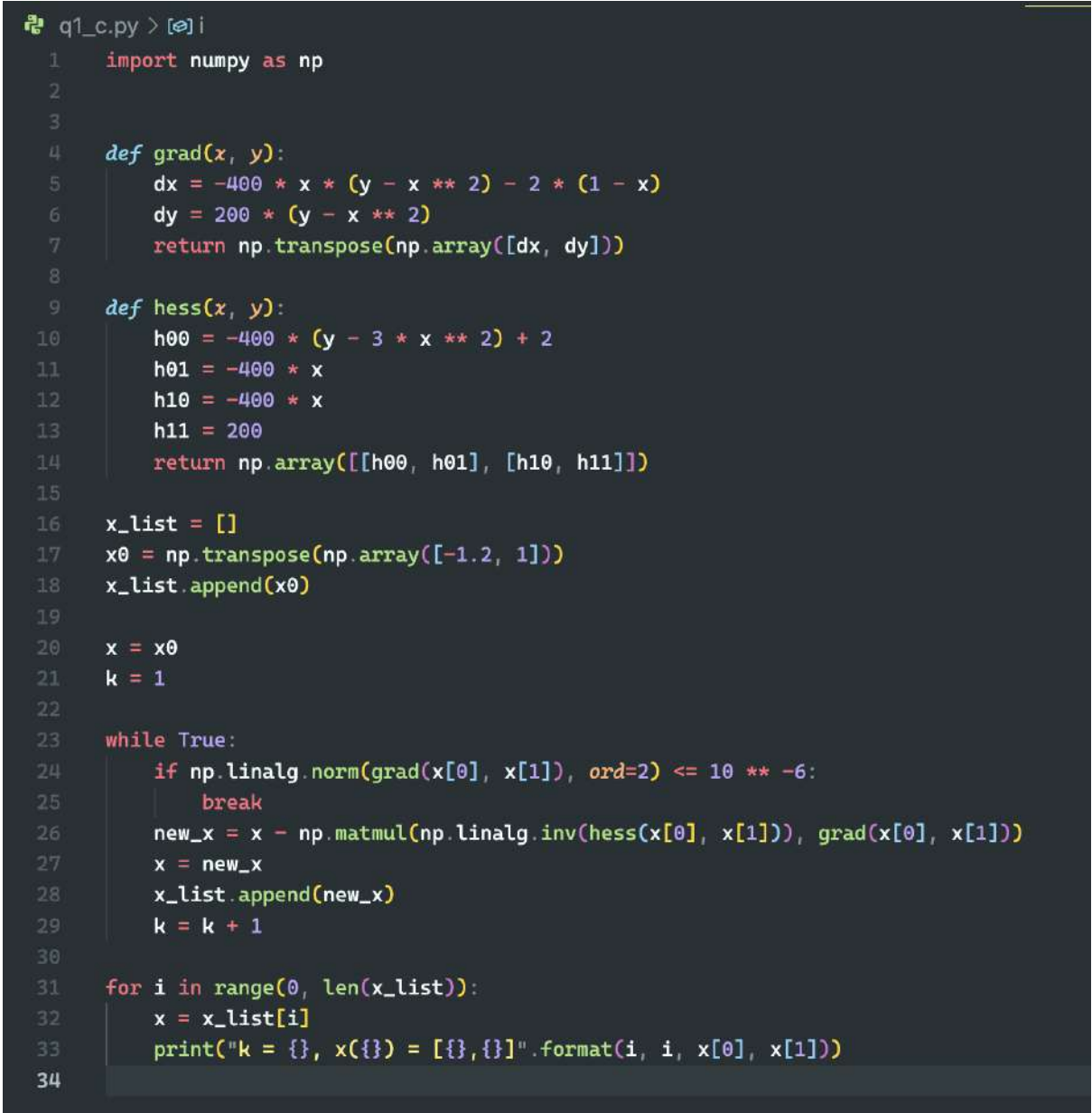
$$\therefore \text{hessian} = \begin{bmatrix} -400(y - 3x^2) + 2 & -400x \\ -400x & 200 \end{bmatrix}$$

c)

The iterations:

```
(base) → hw2 python -u "/Users/yueyifei/Desktop/COMP9417/homework/hw2/q1_c.py"
k = 0, x(0) = [-1.2, 1.0]
k = 1, x(1) = [-1.1752808988764043, 1.3806741573033703]
k = 2, x(2) = [0.763114871176324, -3.175033854747852]
k = 3, x(3) = [0.7634296788839284, 0.5828247754969258]
k = 4, x(4) = [0.999995311085015, 0.9440273238533163]
k = 5, x(5) = [0.999995695653691, 0.9999913913257614]
k = 6, x(6) = [0.9999999999999999, 0.999999999814724]
```

A screen shot of code for c):



```
q1_c.py > [i]
1  import numpy as np
2
3
4  def grad(x, y):
5      dx = -400 * x * (y - x ** 2) - 2 * (1 - x)
6      dy = 200 * (y - x ** 2)
7      return np.transpose(np.array([dx, dy]))
8
9  def hess(x, y):
10     h00 = -400 * (y - 3 * x ** 2) + 2
11     h01 = -400 * x
12     h10 = -400 * x
13     h11 = 200
14     return np.array([[h00, h01], [h10, h11]])
15
16  x_list = []
17  x0 = np.transpose(np.array([-1.2, 1]))
18  x_list.append(x0)
19
20  x = x0
21  k = 1
22
23  while True:
24      if np.linalg.norm(grad(x[0], x[1]), ord=2) <= 10 ** -6:
25          break
26      new_x = x - np.matmul(np.linalg.inv(hess(x[0], x[1])), grad(x[0], x[1]))
27      x = new_x
28      x_list.append(new_x)
29      k = k + 1
30
31  for i in range(0, len(x_list)):
32      x = x_list[i]
33      print("k = {}, x({}) = [{} , {}]".format(i, i, x[0], x[1]))
34
```

Question 2

a)

The coordinate level GD updates:

$$a) L(\beta_0, \beta) = \frac{1}{2} \|\beta\|_2^2 + \frac{\lambda}{n} \sum_{i=1}^n \left[y_i \ln \left(\frac{1}{\sigma(\beta_0 + \beta^T x_i)} \right) + (1 - y_i) \ln \left(\frac{1}{1 - \sigma(\beta_0 + \beta^T x_i)} \right) \right]$$

Let $z_i = \beta_0 + \beta^T x_i$,

Then we have.

$$L(\beta_0, \beta) = \frac{1}{2} \|\beta\|_2^2 - \frac{\lambda}{n} \sum_{i=1}^n \left[y_i \ln \sigma(z_i) + (1 - y_i) \ln (1 - \sigma(z_i)) \right]$$

Therefore, when $j \neq 0$,

$$\begin{aligned} \frac{\partial L(\beta_0, \beta)}{\partial \beta_j} &= \beta_j - \frac{\lambda}{n} \sum_{i=1}^n \left[\frac{y_i}{\sigma(z_i)} - \frac{1 - y_i}{1 - \sigma(z_i)} \right] \sigma(z_i) (1 - \sigma(z_i)) \\ &= \beta_j - \frac{\lambda}{n} \sum_{i=1}^n [y_i - \sigma(z_i)] \\ &= \beta_j - \frac{\lambda}{n} \sum_{i=1}^n \left[y_i - \frac{1}{1 + e^{-\beta_0 - \beta^T x_i}} \right] x_{ij} \\ \therefore \frac{\partial L(\beta_0, \beta)}{\partial \beta_j} &= \beta_j + \frac{\lambda}{n} \sum_{i=1}^n \left[\frac{1}{1 + e^{-\beta_0 - \beta^T x_i}} - y_i \right] x_{ij} \end{aligned}$$

and when $j = 0$,

$$\begin{aligned} \frac{\partial L(\beta_0, \beta)}{\partial \beta_0} &= -\frac{\lambda}{n} \sum_{i=1}^n [y_i - \sigma(z_i)] \\ &= -\frac{\lambda}{n} \sum_{i=1}^n [y_i - \frac{1}{1 + e^{-\beta_0 - \beta^T x_i}}] \\ &= \frac{\lambda}{n} \sum_{i=1}^n \left[\frac{1}{1 + e^{-\beta_0 - \beta^T x_i}} - y_i \right] \end{aligned}$$

$$\therefore \frac{\partial L(\beta_0, \beta)}{\partial \beta_0} = \frac{\lambda}{n} \sum_{i=1}^n \left[\frac{1}{1 + e^{-\beta_0 - \beta^T x_i}} - y_i \right]$$

b)

The vectorized GD updates:

b) Since $L(\beta_0, \beta) = \frac{1}{2} \|\beta\|^2 - \frac{\lambda}{n} \sum_{i=1}^n [y_i \ln \sigma(z_i) + (1-y_i) \ln (1-\sigma(z_i))]$
 $z_i = \beta_0 + \beta^T x_i$

Then, $\frac{\partial L(\beta_0, \beta)}{\partial \beta} = \beta - \frac{\lambda}{n} \sum_{i=1}^n [y_i - \sigma(z_i)]$
 $= \beta + \frac{\lambda}{n} \sum_{i=1}^n \left[\frac{1}{1+e^{\beta_0 + \beta^T x_i}} - y_i \right] x_i^T$
 $= \beta + X^T \left(\frac{1}{1+e^{\beta_0 + \beta^T X}} - y \right)$

$\frac{\partial L(\beta_0, \beta)}{\partial \beta_0} = \frac{\lambda}{n} \sum_{i=1}^n \left[\frac{1}{1+e^{\beta_0 + \beta^T x_i}} - y_i \right]$ (proved in a1)

Therefore, $\begin{bmatrix} \frac{\partial L(\beta_0, \beta)}{\partial \beta_0} & \frac{\partial L(\beta_0, \beta)}{\partial \beta}^T \end{bmatrix}^T$
 $\beta^{(k)} = \beta^{(k-1)} - \gamma \left[\frac{\partial L(\beta_0, \beta)}{\partial \beta_0}, \left(\frac{\partial L(\beta_0, \beta)}{\partial \beta} \right)^T \right]^T$

c)

The vectorized dampened Newton updates:

c) Since $\frac{\partial L(\beta_0, \beta)}{\partial \beta_j} = \beta_j + \frac{\lambda}{n} \sum_{i=1}^n \left[\frac{1}{1 + e^{-\beta_0 - \beta^T x_i}} - y_i \right] x_{ij}$,

let $x_i = [1, x_i]$

then the element $h_{j,k}$ in Hessian matrix is:

$$h_{j,k} = \frac{\partial^2 L(\beta_0, \beta)}{\partial \beta_j \partial \beta_k} = \frac{\lambda}{n} \sum_{i=1}^n \left[\frac{x_{ik} x_{ij} \cdot e^{-\beta_0 - \beta^T x_i}}{(1 + e^{-\beta_0 - \beta^T x_i})^2} \right] \quad (j \neq k)$$

and

$$h_{j,j} = \frac{\partial^2 L(\beta_0, \beta)}{\partial \beta_j \partial \beta_j} = 1 + \frac{\lambda}{n} \sum_{i=1}^n \left[\frac{x_{ij} x_{ij} \cdot e^{-\beta_0 - \beta^T x_i}}{(1 + e^{-\beta_0 - \beta^T x_i})^2} \right] \quad (j=k)$$

(elements on diagonal)

~~Then we have.~~

Therefore, $H = I + X^T S X$

where $I = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{bmatrix}$

$$S = \frac{\lambda}{n} \begin{bmatrix} \sigma(z_1)(1-\sigma(z_1)) & \dots & 0 \\ 0 & \dots & \sigma(z_2)(1-\sigma(z_2)) & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma(z_n)(1-\sigma(z_n)) \end{bmatrix}$$

$$\therefore \gamma^{(k)} = \gamma^{(k-1)} - \partial X H^{-1} \left[\frac{\partial L}{\partial \beta_0}, \left(\frac{\partial L}{\partial \beta} \right)^T \right]^T$$

d)

A print out of the rows requested:

```
(base) → hw2 python -u "/Users/yueyifei/Desktop/COMP9417/homework/hw2/q2_def.py"
The first row of X_train:
[-0.93555843  0.67519298  1.3849985 ]
The last row of X_train:
[-1.13301479 -1.09458877  0.96702449]

The first row of X_test:
[-0.29382524  1.36005105  0.26306826]
The last row of X_test:
[-0.29382524 -1.05390413 -1.34833155]

The first row of Y_train:
0
The last row of Y_train:
1

The first row of Y_test:
0
The last row of Y_test:
1
```

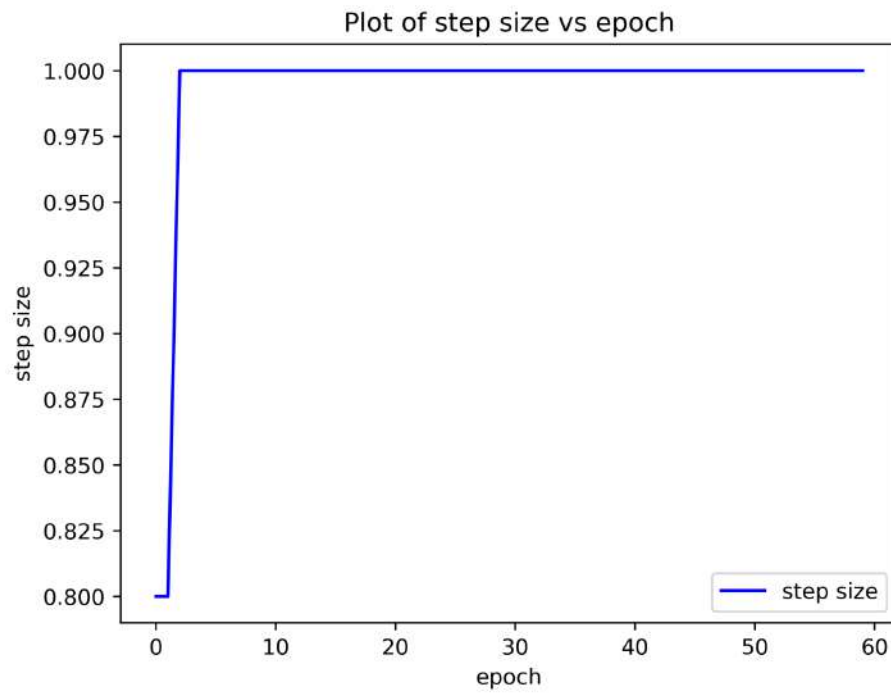
A screen shot of code for d):


```
q2_def.py > ...
1 import numpy as np
2 import pandas as pd
3 from sklearn.model_selection import train_test_split
4 from sklearn.preprocessing import StandardScaler
5 from helper import sigmoid, loss
6 import matplotlib.pyplot as plt
7
8 # Question d of Part 2
9 data = pd.read_csv('songs.csv')
10
11 # (I)
12 remove_features = ['Artist Name', 'Track Name', 'key', 'mode', 'time_signature', 'instrumentalness']
13 data = data.drop(remove_features, axis=1)
14
15 # (II)
16 data = data[(data['Class'] == 5) | (data['Class'] == 9)]
17 data['Class'].replace([5, 9], [1, 0], inplace=True)
18 data = data.reset_index(drop=True)
19
20 # (III)
21 data = data.dropna()
22
23 # (IV)
24 X_train, X_test, Y_train, Y_test = train_test_split(data.iloc[:, 0:-1], data.iloc[:, -1], test_size=0.3,
25                                                    random_state=23)
26
27 # (V)
28 scaler = StandardScaler()
29 scaler.fit(X_train)
30 X_train = scaler.transform(X_train)
```

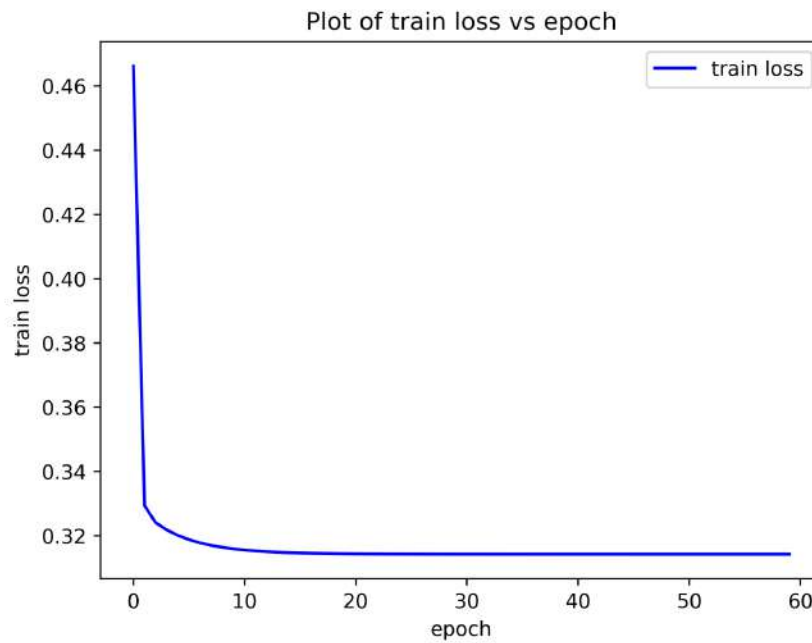
```
31
32     # (VI)
33     Y_train = Y_train.to_numpy()
34     Y_test = Y_test.to_numpy()
35
36     print('The first row of X_train:')
37     print(X_train[0][:3])
38     print('The last row of X_train:')
39     print(X_train[-1][:3])
40     print()
41     print('The first row of X_test:')
42     print(X_test[0][:3])
43     print('The last row of X_test:')
44     print(X_test[-1][:3])
45     print()
46     print('The first row of Y_train:')
47     print(Y_train[0])
48     print('The last row of Y_train:')
49     print(Y_train[-1])
50     print()
51     print('The first row of Y_test:')
52     print(Y_test[0])
53     print('The last row of Y_test:')
54     print(Y_test[-1])
55     print()
56
```

e)

The plot of step size vs epoch:



The plot of train loss vs epoch:



Final achieved train and test losses:

```
Final loss achieved by GD algorithm on the train data
0.3142968229209097
Final loss achieved by GD algorithm on the test data
0.31721131314546464
```

A screen shot of code for e):

```
58 # Question e of Part 2
59
60 epochs = 60
61 lam = 0.5
62 beta_0 = 0
63 beta = np.ones(len(data.columns) - 1)
64 gamma_0 = np.insert(beta, 0, beta_0)
65 a = 0.5
66 b = 0.8
67
68
69 def grad(gamma, X, y, lam):
70     z = np.dot(X, gamma[1:]) + gamma[0]
71     grad_beta = gamma[1:] + (lam / len(X)) * np.dot((sigmoid(z) - y), X)
72     grad_beta0 = (lam / len(X)) * np.sum((sigmoid(z) - y))
73     grad_gamma = np.insert(grad_beta, 0, grad_beta0)
74     return grad_gamma
75
76
77 def get_gd_gammas_alphas(gamma_0, X, y, epochs, lam):
78     alphas = []
79     gammas = []
80     gamma = gamma_0
81     alphas.append(1)
82     gammas.append(gamma)
83     # Backtracking line search
84     for k in range(epochs):
85         alpha = 1
86         if loss(gamma - alpha * grad(gamma, X, y, lam), X, y, lam) > loss(gamma, X, y, lam) - a * alpha * np.linalg.norm(grad(gamma, X, y, lam), ord=2)**2:
87             alpha = alpha * b
88             alphas.append(alpha)
89             gamma_updated = gamma - alpha * grad(gamma, X, y, lam)
90             gamma = gamma_updated
91             gammas.append(gamma_updated)
92     return gammas, alphas
93
```

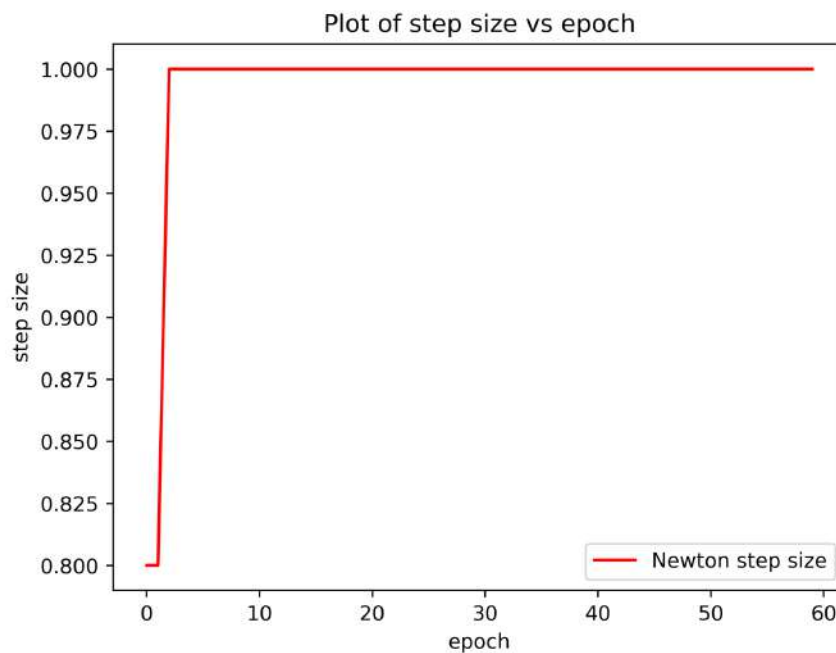
```

94
95 def get_gd_losses(gamma_0, X, y, epochs, lam):
96     losses = []
97     gammas = get_gd_gammas_alphas(gamma_0, X, y, epochs, lam)[0]
98     for gamma in gammas:
99         gamma_loss = loss(gamma, X, y, lam)
100         losses.append(gamma_loss)
101     return losses
102
103
104 alphas = get_gd_gammas_alphas(gamma_0, X_train, Y_train, epochs, lam)[1]
105 losses_train = get_gd_losses(gamma_0, X_train, Y_train, epochs, lam)
106 losses_test = get_gd_losses(gamma_0, X_test, Y_test, epochs, lam)
107 print('Final loss achieved by GD algorithm on the train data')
108 print(losses_train[-1])
109 print('Final loss achieved by GD algorithm on the test data')
110 print(losses_test[-1])
111
112 plt.plot(range(epochs), alphas[1:], color='blue', label='step size')
113 plt.xlabel('epoch')
114 plt.legend() # creates legend in top right corner of plot
115 plt.ylabel('step size')
116 plt.title('Plot of step size vs epoch')
117 plt.savefig('qe_step_sizes.png', dpi=400)
118 plt.show()
119
120 plt.plot(range(epochs), losses_train[1:], color='blue', label='train loss')
121 plt.xlabel('epoch')
122 plt.legend() # creates legend in top right corner of plot
123 plt.ylabel('train loss')
124 plt.title('Plot of train loss vs epoch')
125 plt.savefig('qe_train_loss.png', dpi=400)
126 plt.show()
127

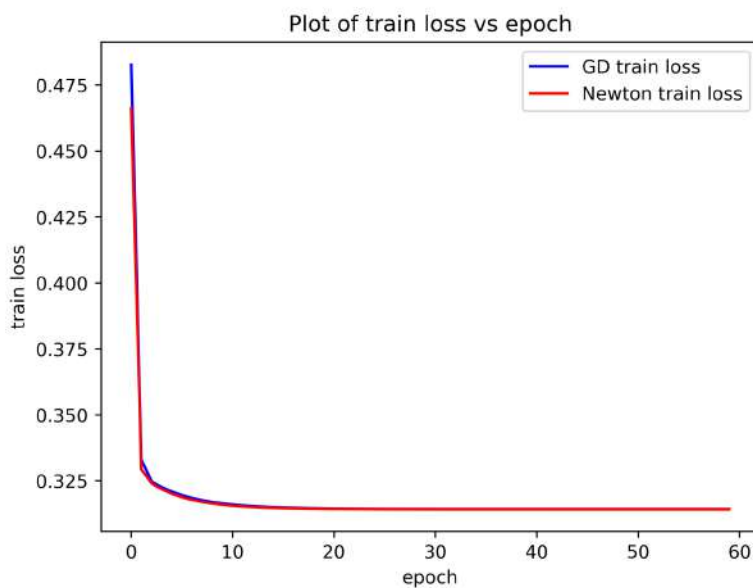
```

f)

The plot of step size vs epoch:



The plot of train losses (of both GD and Newton) vs epoch:



Final achieved train and test losses:

Final loss achieved by Newton algorithm on the train data
0.31429685408159996
Final loss achieved by Newton algorithm on the test data
0.3172113358850347

A screen shot of code for f):

```
60 epochs = 60
61 lam = 0.5
62 beta_0 = 0
63 beta = np.ones(len(data.columns) - 1)
64 gamma_0 = np.insert(beta, 0, beta_0)
65 a = 0.5
66 b = 0.8
67
68
69 def grad(gamma, X, y, lam):
70     z = np.dot(X, gamma[1:]) + gamma[0]
71     grad_beta = gamma[1:] + (lam / len(X)) * np.dot((sigmoid(z) - y), X)
72     grad_beta0 = (lam / len(X)) * np.sum((sigmoid(z) - y))
73     grad_gamma = np.insert(grad_beta, 0, grad_beta0)
74     return grad_gamma
75
```

```
129 # Question f of Part 2
130
131 def hess(gamma, X, lam):
132     # H = I + X.T * S * X
133     hess_matrix = np.zeros((11, 11))
134     identity_matrix = np.identity(11)
135     z = np.dot(X, gamma[1:]) + gamma[0]
136     diagnose_matrix = np.diag((lam / len(X)) * sigmoid(z) * (1 - sigmoid(z)))
137     col = np.ones((len(X), 1))
138     X = np.column_stack((col, X))
139     hess_matrix = identity_matrix + X.T @ diagnose_matrix @ X
140     return hess_matrix
141
142
143 def get_dn_gammas_alphas(gamma_0, X, y, epochs, lam):
144     alphas = []
145     gammas = []
146     gamma = gamma_0
147     alphas.append(1)
148     gammas.append(gamma)
149     # Backtracking line search
150     for k in range(epochs):
151         alpha = 1
152         if loss(gamma - alpha * grad(gamma, X, y, lam), X, y, lam) > loss(gamma, X, y, lam) - a * alpha * np.linalg.norm(grad(gamma, X, y, lam), ord=2)**2:
153             alpha = alpha * b
154         alphas.append(alpha)
155         gamma_updated = gamma - alpha * np.matmul(np.linalg.inv(hess(gamma, X, lam)), grad(gamma, X, y, lam))
156         gamma = gamma_updated
157         gammas.append(gamma_updated)
158     return gammas, alphas
159
```

```

160
161 def get_dn_losses(gamma_0, X, y, epochs, lam):
162     losses = []
163     gammas = get_dn_gammas_alphas(gamma_0, X, y, epochs, lam)[0]
164     for gamma in gammas:
165         gamma_loss = loss(gamma, X, y, lam)
166         losses.append(gamma_loss)
167     return losses
168
169
170 alphas = get_dn_gammas_alphas(gamma_0, X_train, Y_train, epochs, lam)[1]
171 losses_train_dn = get_dn_losses(gamma_0, X_train, Y_train, epochs, lam)
172 losses_test_dn = get_dn_losses(gamma_0, X_test, Y_test, epochs, lam)
173 print('Final loss achieved by Newton algorithm on the train data')
174 print(losses_train_dn[-1])
175 print('Final loss achieved by Newton algorithm on the test data')
176 print(losses_test_dn[-1])
177
178 plt.plot(range(epochs), alphas[1:], color='red', label='Newton step size')
179 plt.xlabel('epoch')
180 plt.legend() # creates legend in top right corner of plot
181 plt.ylabel('step size')
182 plt.title('Plot of step size vs epoch')
183 plt.savefig('qf_step_sizes.png', dpi=400)
184 plt.show()
185
186 plt.plot(range(epochs), losses_train_dn[1:], color='blue', label='GD train loss')
187 plt.plot(range(epochs), losses_train[1:], color='red', label='Newton train loss')
188 plt.xlabel('epoch')
189 plt.legend() # creates legend in top right corner of plot
190 plt.ylabel('train loss')
191 plt.title('Plot of train loss vs epoch')
192 plt.savefig('qf_train_loss.png', dpi=400)
193 plt.show()
194

```

g)

Because Newton's method requires to calculate the Hessian matrix, but the time complexity of calculating Hessian matrix is higher. It is often more difficult or intractable to derive the second derivative, which requires a lot of computation. If n values are required to compute the first derivative, then n^2 are required for the second derivative. Therefore, GD and its variants are much more popular than Newton's method for solving machine learning problems.

Question 3

a)

a) For $\hat{\mu}_{MLE}$

since $\hat{\mu}_{MLE} = \bar{X}$,

$\text{bias}(\hat{\mu}_{MLE}) = \text{bias}(\bar{X})$

$$= E\left(\frac{1}{n} \sum_{i=1}^n X_i\right) - \mu$$

$$= \frac{1}{n} \sum_{i=1}^n E(X_i) - \mu$$

$$= \frac{1}{n} \sum_{i=1}^n \mu - \mu$$

$$= \mu - \mu$$

$$= 0$$

say \bar{X} is an unbiased estimator for μ

$\text{var}(\hat{\mu}_{MLE}) = \text{var}(\bar{X})$

$$= \text{var}\left(\frac{1}{n} \sum_{i=1}^n X_i\right)$$

$$= \frac{1}{n^2} \sum_{i=1}^n \text{var}(X_i)$$

$$= \frac{1}{n}$$

For $\hat{\sigma}_{MLE}^2$

since $\hat{\sigma}_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$

$\text{bias}(\hat{\sigma}_{MLE}^2) = \text{bias}\left(\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2\right)$

$$= E\left(\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2\right) - \sigma^2$$

$$= E\left(\frac{1}{n} \sum_{i=1}^n (X_i - \mu + \mu - \bar{X})^2\right) - \sigma^2$$

$$= \frac{1}{n} E\left(\sum_{i=1}^n ((X_i - \mu)^2 - 2(X_i - \mu)(\bar{X} - \mu) + (\bar{X} - \mu)^2)\right) - \sigma^2$$

$$= \frac{1}{n} E\left(\sum_{i=1}^n (X_i - \mu)^2 - 2n(\bar{X} - \mu)(\bar{X} - \mu) + n(\bar{X} - \mu)^2\right) - \sigma^2$$

$$= \frac{1}{n} E\left(\sum_{i=1}^n (X_i - \mu)^2 - n(\bar{X} - \mu)^2\right) - \sigma^2$$

$$= \frac{1}{n} \left(\sum_{i=1}^n E((X_i - \mu)^2) - nE((\bar{X} - \mu)^2)\right) - \sigma^2$$

$$= \frac{1}{n} (n \text{Var}(X) - n \text{Var}(\bar{X})) - \sigma^2$$

$$= \text{Var}(X) - \text{Var}(\bar{X}) - \sigma^2$$

$$= \sigma^2 - \frac{\sigma^2}{n} - \sigma^2$$

$$= -\frac{\sigma^2}{n}$$

For $\text{Var}(\hat{\sigma}_{MLE}^2) = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2\right)$

since $\text{Var}\left(\frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \bar{X})^2\right) = 2(n-1)$

$$\text{Var}\left(\sum_{i=1}^n (X_i - \bar{X})^2\right) = 2\sigma^4(n-1)$$

$$\text{Var}\left(\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2\right) = \frac{2\sigma^4(n-1)}{n^2}$$

so $\text{Var}(\hat{\sigma}_{MLE}^2) = \frac{2\sigma^4(n-1)}{n^2}$

b)

The bias and variance of the new estimator:

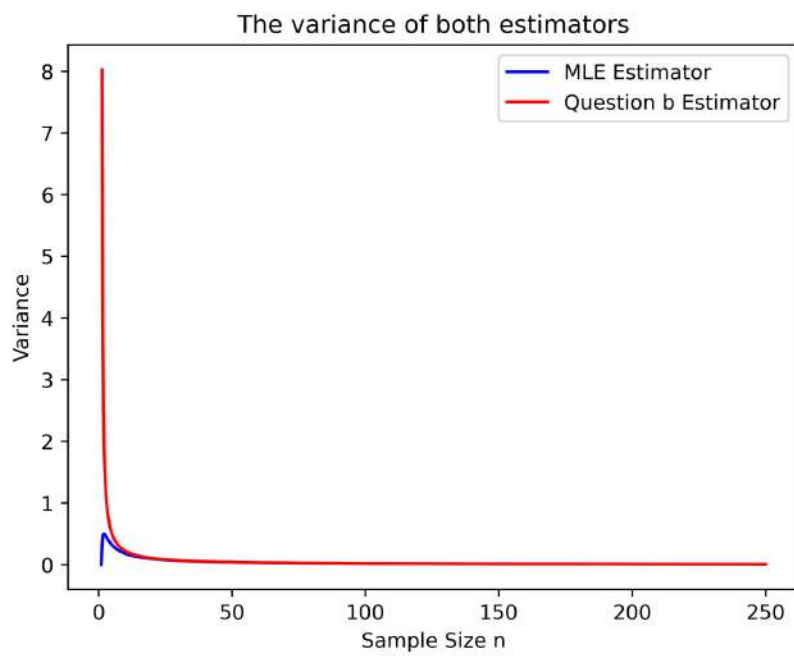
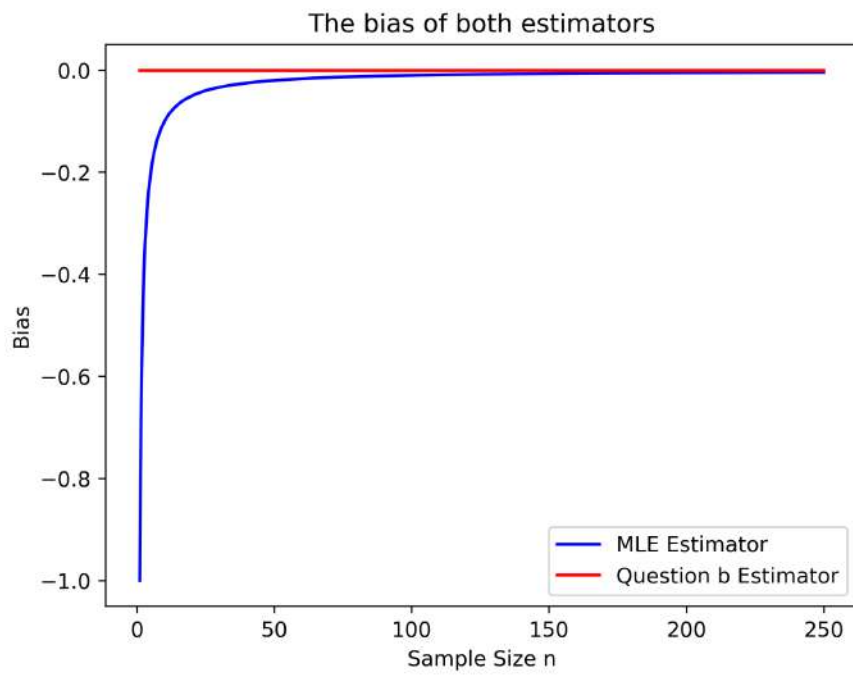
$$\begin{aligned}
 b) \quad \text{bias}(\bar{\sigma}^2) &= \text{bias}\left(\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2\right) \\
 &= E\left(\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2\right) - \sigma^2 \\
 &= \frac{1}{n-1} E\left(\sum_{i=1}^n (X_i - \mu)^2 - n(\bar{X} - \mu)^2\right) - \sigma^2 \quad (\text{as we proved in a}) \\
 &= \frac{1}{n-1} (n \text{Var}(X) - n \text{Var}(\bar{X})) - \sigma^2 \\
 &= \frac{n \text{Var}(X)}{n-1} - \frac{n \text{Var}(\bar{X})}{n-1} - \sigma^2 \\
 &= \frac{n \sigma^2}{n-1} - \frac{\sigma^2}{n-1} - \sigma^2 \\
 &= 0
 \end{aligned}$$

$$\text{Var}(\bar{\sigma}^2) = \text{Var}\left(\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2\right)$$

$$\text{Since } \text{Var}\left(\frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \bar{X})^2\right) = 2(n-1)$$

$$\text{Var}\left(\sum_{i=1}^n (X_i - \bar{X})^2\right) = 2\sigma^4(n-1)$$

$$\text{Var}\left(\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2\right) = \frac{2\sigma^4}{n-1}$$



A screen shot of code for b):

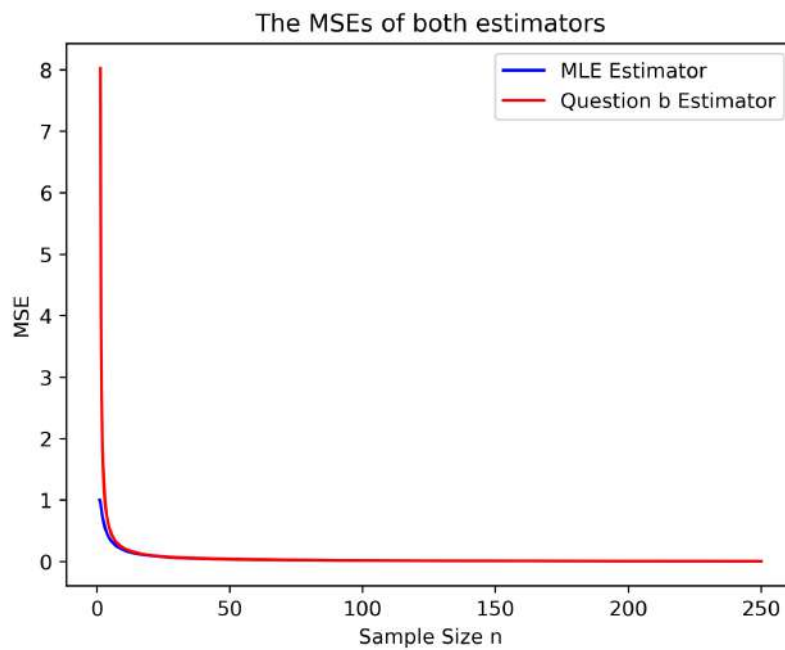
```
q3_bc.py > [MSE_BE]
1  import numpy as np
2  import matplotlib.pyplot as plt
3
4  # Question b of Part 3
5  sigma = 1
6
7  bias_MLE = lambda n: (-1 * sigma ** 2) / n
8  var_MLE = lambda n: (2 * sigma ** 4 * (n - 1)) / n ** 2
9
10 bias_BE = lambda n: n * 0
11 var_BE = lambda n: (2 * sigma ** 4) / (n - 1)
12
13 nrange = np.linspace(1, 250, 1000)
14
15 # Plot the bias of both estimators
16 plt.plot(nrange, bias_MLE(nrange), label="MLE Estimator", color='blue')
17 plt.plot(nrange, bias_BE(nrange), label="Question b Estimator", color="red")
18
19 plt.xlabel('Sample Size n')
20 plt.ylabel('Bias')
21 plt.title('The bias of both estimators')
22 plt.legend()
23 plt.savefig('q3b_bias.png', dpi=400)
24 plt.show()
25
26 # the variance of both estimators
27 plt.plot(nrange, var_MLE(nrange), label="MLE Estimator", color='blue')
28 plt.plot(nrange, var_BE(nrange), label="Question b Estimator", color="red")
29
30 plt.xlabel('Sample Size n')
31 plt.ylabel('Variance')
32 plt.title('The variance of both estimators')
33 plt.legend()
34 plt.savefig('q3b_var.png', dpi=400)
35 plt.show()
36
```


c)

$$\begin{aligned} \text{c) MLE: } \text{MSE}(\hat{\sigma}^2) &= \left(-\frac{\sigma^2}{n}\right)^2 + \frac{2\sigma^4(n-1)}{n^2} \\ &= \frac{\cancel{2}\sigma^4}{\cancel{n^2}} \frac{\sigma^4(2n-1)}{n^2} \end{aligned}$$

~~Old~~ New Estimator:

$$\begin{aligned} \text{MSE}(\hat{\sigma}^2) &= \sigma^2 + \frac{2\sigma^4}{n-1} \\ &= \frac{2\sigma^4}{n-1} \end{aligned}$$



A screen shot of code for c):

```
37 # Question c of Part 3
38 MSE_MLE = lambda n: bias_MLE(n) ** 2 + var_MLE(n)
39 MSE_BE = lambda n: bias_BE(n) ** 2 + var_BE(n)
40
41 # the MSEs of both estimators
42 plt.plot(nrange, MSE_MLE(nrange), label="MLE Estimator", color='blue')
43 plt.plot(nrange, MSE_BE(nrange), label="Question b Estimator", color="red")
44
45 plt.xlabel('Sample Size n')
46 plt.ylabel('MSE')
47 plt.title('The MSEs of both estimators')
48 plt.legend()
49 plt.savefig('q3c.png', dpi=400)
50 plt.show()
51
```

Discussion:

I think MLE is better. With the increase of the size of samples, the gap between the two estimators will be smaller and smaller, but when the size of samples is small, the MSE of MLE is smaller, so MLE is better to some extent.