

① a) we can solve that using normal equations:

$$A = X^T X = \begin{pmatrix} 1.5 & 3 & 4.5 \\ 2 & 2.5 & 3 \\ 4.5 & 3 & 3 \end{pmatrix} \begin{pmatrix} 1.5 & 2 \\ 3 & 2.5 \\ 4.5 & 3 \end{pmatrix} = \begin{pmatrix} 31.5 & 24 \\ 24 & 17.25 \end{pmatrix}$$

$$B = X^T Y = \begin{pmatrix} 1.5 & 3 & 4.5 \\ 2 & 2.5 & 3 \end{pmatrix} \begin{pmatrix} 10 \\ 15.5 \\ 21 \end{pmatrix} = \begin{pmatrix} 156 \\ 121.75 \end{pmatrix}$$

$$\text{Solve } Ax = B \rightarrow x = A^{-1} B = \begin{pmatrix} 2.6 \\ 3 \end{pmatrix} = \beta$$

② First of all, in the Part a), the analytical solution seems easy to find. But it is not the case for all ML problems. Some cases the analytical solution might be very complex to find. The second reason can be that in such cases, these kinds of solutions might be computationally infeasible for current available compute powers. But learning algorithms can be implemented without this issue. Moreover, they can also be designed to be flexible.

$$\textcircled{c} L = \frac{1}{2} \|X\beta - y\|_2^2 \rightarrow \nabla L(\beta) = 2(X^T X \beta - X^T y)$$

$$\mu = 0.1, \beta_0 = (1, 1)^T \rightarrow \text{error} = L(\beta_0) = 108.16$$

$$\nabla L = \begin{pmatrix} -201 \\ -157 \end{pmatrix}$$

$$\beta_1 = \beta_0 - \mu \nabla L = \begin{pmatrix} 21.1 \\ 16.7 \end{pmatrix} \rightarrow \text{error} = L(\beta_1) = 8812.7$$

$$\nabla L = \begin{pmatrix} 1818.9 \\ 1412.25 \end{pmatrix}$$

$$\beta_2 = \beta_1 - \mu \nabla L = \begin{pmatrix} -160.79 \\ -124.52 \end{pmatrix} \rightarrow \text{error} = 718408.69 = L(\beta_2)$$

```
1 # the code used for gredient descent
2
3 import numpy as np
4 b = np.array([[1], [1]], dtype='float')
5 x = np.array([[1.5,2], [3,2.5], [4.5,3]])
6 y = np.array([[10],[15.5],[21]])
7 moo = 0.1
8 def grad(x,y,b):
9     return 2*(x.T@x@b-x.T@y)
10
11 def error(x, y , b):
12     return np.mean((x@b - y) ** 2)
13
14 print(f'b : {b} | error : {error(x, y , b)}')
15 b1 = b - moo * grad(x,y,b)
16
17 print(f'b1 : {b1} | error : {error(x, y , b1)}')
18 b2 = b1 - moo * grad(x,y,b1)
19 print(f'b2 : {b2} | error : {error(x, y , b2)}')
```

②
a) In stochastic gradient descent (SGD), we don't go through the whole dataset in one update step. In other words, the loss function would not have a sum over N , but rather over a mini-batch of 6, 32, 64, 128 data points. That's because with big datasets and complex models, computing gradients in that way would be computationally expensive and even not feasible.

b) The code and intermediate results for this part are given in the next page.

c) The code and intermediate results are given in the next pages.

→ Also Adagrad helps in this algorithm. We can see that by comparing losses.

The code and the results (final and intermediate) for question 2, part b)

```
1  # the code used for stochastic gradient descent
2
3  import numpy as np
4  b = np.array([[1], [1]], dtype='float')
5  x = np.array([[1.5,2], [3,2.5], [4.5,3]])
6  y = np.array([[10],[15.5],[21]])
7  moo = 0.1
8  epochs = 2
9
10 def grad(x, y,b):
11     return 2*(x.T@x@b-x.T@y)
12
13 def mse(x,y,b):
14     return (x@b - y)**2
15
16 for epoch in range(epochs):
17     for i in range(len(x)):
18         x_i = x[i, :].reshape((1,2))
19         y_i = y[i, :].reshape((1,1))
20         gradient = -grad(x_i, y_i,b)
21         err = mse(x_i, y_i,b)
22         print(f'epoch {epoch} - step {i}\n=====')
23         print(f'b: {b}')
24         print(f'mse: {err}')
25         b = b - moo * gradient
26
27 print(f'final b is : {b}')
```

```
epoch 0 - step 0
=====
b: [[1.]
     [1.]]
mse: [[42.25]]
epoch 0 - step 1
=====
b: [[-0.95]
     [-1.6 ]]
mse: [[499.5225]]
epoch 0 - step 2
=====
b: [[-14.36 ]
     [-12.775]]
mse: [[15362.363025]]
epoch 1 - step 0
=====
b: [[-125.9105]
     [ -87.142 ]]
mse: [[139240.73592506]]
epoch 1 - step 1
=====
b: [[-237.855425]
     [-236.4019  ]]
mse: [[1742587.51104455]]
epoch 1 - step 2
=====
b: [[-1029.89804  ]
     [ -896.4374125]]
mse: [[53946871.72456145]]
final b is : [[-7640.26611575]
               [-5303.349463  ]]
```

The code and the results (final and intermediate) for question 2, part c)

```
1  # the code used for stochastic gradient descent with adagrad
2  import numpy as np
3  b = np.array([[1], [1]], dtype='float')
4  x = np.array([[1.5,2], [3,2.5], [4.5,3]])
5  y = np.array([[10],[15.5],[21]])
6  moo = 0.1
7  epochs = 2
8  epsilon = 1e-8
9  def grad(x, y,b):
10     return 2*(x.T@x@b-x.T@y)
11
12  def mse(x,y,b):
13     return (x@b - y)**2
14  gradian_history = np.array([[0],[0]], dtype='float')
15  for epoch in range(epochs):
16     for i in range(len(x)):
17         x_i = x[i, :].reshape((1,2))
18         y_i = y[i, :].reshape((1,1))
19         gradient = -grad(x_i, y_i,b)
20         err = mse(x_i, y_i,b)
21         print(f'epoch {epoch} - step {i}\n=====')
22         print(f'b: {b}')
23         print(f'mse: {err}')
24         gradian_history += gradient**2
25         step_size = moo / (np.sqrt(gradian_history) + epsilon)
26         b = b - step_size * gradient
27  print(f'final b is : {b}')
```

```
epoch 0 - step 0
=====
b: [[1.]
     [1.]]
mse: [[42.25]]
epoch 0 - step 1
=====
b: [[0.9]
     [0.9]]
mse: [[111.30249999]]
epoch 0 - step 2
=====
b: [[0.8044319 ]
     [0.81030368]]
mse: [[223.47694926]]
epoch 1 - step 0
=====
b: [[0.71471464]
     [0.72667635]]
mse: [[55.86927667]]
epoch 1 - step 1
=====
b: [[0.69992617]
     [0.69982432]]
mse: [[135.7378943]]
epoch 1 - step 2
=====
b: [[0.65805937]
     [0.65346741]]
mse: [[258.51271487]]
final b is : [[0.59256638]
               [0.59257312]]
```