
CS 189: INTRODUCTION TO
MACHINE LEARNING

Fall 2017



HOMEWORK 6



Solutions by

JINHONG DU

3033483677

Question 1

(a)

Jinhong Du
jaydu@berkeley.edu

(b)

I certify that all solutions are entirely in my words and that I have not looked at another student's solutions. I have credited all external sources in this write up.

Jinhong Du

Question 2

(a)

$$\because \forall x_1, x_2 \in \mathbb{R}^d, \forall \lambda \in \mathbb{R}, \lambda \in [0, 1],$$

$$\begin{aligned} f[\lambda x_1 + (1 - \lambda)x_2] &= \|\lambda x_1 + (1 - \lambda)x_2 - b\|_2 \\ &= \|\lambda(x_1 - b) + (1 - \lambda)(x_2 - b)\|_2 \\ &\leq \|\lambda x_1\|_2 + \|(1 - \lambda)x_2 - b\|_2 \\ &= |\lambda| \|x_1 + (1 - \lambda)x_2 - b\|_2 + |1 - \lambda| \|x_2 - b\|_2 \\ &= \lambda f(x_1) + (1 - \lambda)f(x_2) \end{aligned}$$

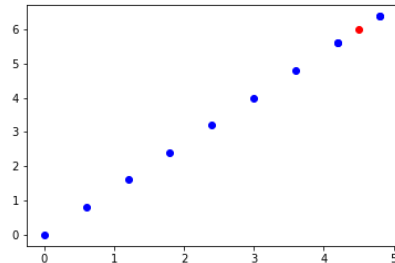
$\therefore f(x)$ is a convex function of x

(b)

$$\begin{aligned} f(x) &= \|x - b\|_2 \\ &= \sqrt{(x_1 - 4.5)^2 + (x_2 - 6)^2} \geq 0 \end{aligned}$$

the equation holds when $x^* = (4.5 \ 6)^T$

```
[ [ 0.  0. ]
 [ 0.6 0.8]
 [ 1.2 1.6]
 [ 1.8 2.4]
 [ 2.4 3.2]
 [ 3.  4. ]
 [ 3.6 4.8]
 [ 4.2 5.6]
 [ 4.8 6.4]
 [ 4.2 5.6]
 [ 4.8 6.4] ]
```



It didn't find the optimal solution. It is because that a step size is too big and x_i moves left and right of the optimal point $(4.5 \ 6)^T$. At the end, x_i change from either $(4.2 \ 5.6)^T$ or $(4.8 \ 6.4)^T$ to the other.

To general $b \neq \vec{0}$, it depends on the choice of b . For example, if $b = (0.6 \ 0.8)^T$, then we can get the optimal solution at the first step. However, if we choose $b = (4.5 \ 6)^T$, it will fail.

\therefore

$$\nabla f(x) = \frac{x - b}{\|x - b\|_2}$$

\therefore

$$x_{i+1} = x_i - \frac{x_i - b}{\|x_i - b\|_2}$$

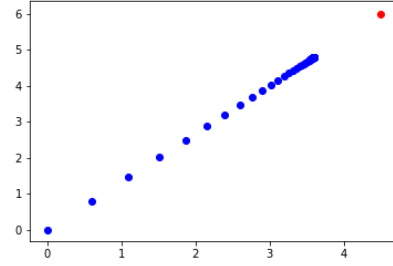
Solution (cont.)

We know that $\frac{x_i - b}{\|x_i - b\|_2}$ is a unit vector with the direction $\pm \overrightarrow{Ox^*} = \pm(4.5, 6) = \pm(0.6, 0.8)$. And $4.5 = 0.6 \times 7 + 0.3 = 0.6 \times 8 - 0.3$, $6.0 = 0.8 \times 7 + 0.4 = 0.8 \times 8 - 0.4$, so the optimal solution will switch between $(4.2, 5.6)$ and $(4.8, 6.4)$.

(c)

After running for 10000 steps, the optimal solution of gradient descend approximates the point $\begin{pmatrix} 3.6 & 4.8 \end{pmatrix}^T$.

```
[ [ 0.      0.      ]
 [ 0.6     0.8     ]
 [ 1.1     1.46666667 ]
 ...
 [ 3.6     4.8     ]
 [ 3.6     4.8     ]
 [ 3.6     4.8     ]]
```



The gradient descend cannot find the optimal solution because the step size is getting smaller and smaller such that it moves slowly from x_i to x_{i+1} when i is large.

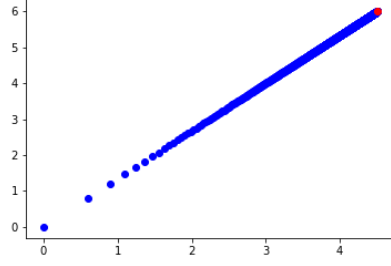
Suppose that $\|x_i\|_2 < \left\| \begin{pmatrix} 4.5 \\ 6 \end{pmatrix} \right\|_2$,

$$\begin{aligned} x_{i+1} &= x_i - \left(\frac{5}{6} \right)^i \frac{x_i - b}{\|x_i - b\|_2} \\ &= x_i + \left(\frac{5}{6} \right)^i \begin{pmatrix} 0.6 \\ 0.8 \end{pmatrix} \\ &= x_0 + \sum_{j=0}^i \left(\frac{5}{6} \right)^j \begin{pmatrix} 0.6 \\ 0.8 \end{pmatrix} \\ \lim_{i \rightarrow \infty} x_{i+1} &= \sum_{j=0}^{\infty} \left(\frac{5}{6} \right)^j \begin{pmatrix} 0.6 \\ 0.8 \end{pmatrix} \\ &= \begin{pmatrix} 3.6 \\ 4.8 \end{pmatrix} \end{aligned}$$

$$\left\| \begin{pmatrix} 3.6 \\ 4.8 \end{pmatrix} \right\|_2 < \left\| \begin{pmatrix} 4.5 \\ 6 \end{pmatrix} \right\|_2$$

so the assumption holds.

(d)



The gradient descend find the optimal solution. At x_{1005} it will get within 0.01 of the optimal solution.

Suppose that $\|x_i\|_2 < \left\| \begin{pmatrix} 4.5 \\ 6 \end{pmatrix} \right\|_2$,

$$\begin{aligned}
 x_{i+1} &= x_i - \left(\frac{1}{i+1} \right)^i \frac{x_i - b}{\|x_i - b\|_2} \\
 &= x_i + \left(\frac{1}{i+1} \right)^i \begin{pmatrix} 0.6 \\ 0.8 \end{pmatrix} \\
 &= x_0 + \sum_{j=0}^i \left(\frac{1}{j+1} \right)^j \begin{pmatrix} 0.6 \\ 0.8 \end{pmatrix} \\
 \lim_{i \rightarrow \infty} x_{i+1} &= \sum_{j=0}^{\infty} \left(\frac{1}{j+1} \right)^j \begin{pmatrix} 0.6 \\ 0.8 \end{pmatrix} \\
 &= \begin{pmatrix} \infty \\ \infty \end{pmatrix} \\
 \left\| \begin{pmatrix} \infty \\ \infty \end{pmatrix} \right\|_2 &> \left\| \begin{pmatrix} 4.5 \\ 6 \end{pmatrix} \right\|_2
 \end{aligned}$$

so the assumption don't holds. So $\exists n_0 \in \mathbb{N}$, s.t. $\left\| \sum_{j=0}^{n_0} \left(\frac{1}{j+1} \right)^j \begin{pmatrix} 0.6 \\ 0.8 \end{pmatrix} \right\|_2 \leq \left\| \begin{pmatrix} 4.5 \\ 6 \end{pmatrix} \right\|_2$ and $\left\| \sum_{j=0}^{n_0+1} \left(\frac{1}{j+1} \right)^j \begin{pmatrix} 0.6 \\ 0.8 \end{pmatrix} \right\|_2 > \left\| \begin{pmatrix} 4.5 \\ 6 \end{pmatrix} \right\|_2$. If the equation holds, then we find the optimal solution at the n_0 th step. Otherwise, $\exists n_1 > n_0$, s.t.

$$\left\| \left[\sum_{j=0}^{n_0} \left(\frac{1}{j+1} \right)^j - \sum_{j=n_0+1}^{n_1} \left(\frac{1}{j+1} \right)^j + \dots \right] \begin{pmatrix} 0.6 \\ 0.8 \end{pmatrix} \right\|_2 \geq \left\| \begin{pmatrix} 4.5 \\ 6 \end{pmatrix} \right\|_2$$

and

$$\left\| \left[\sum_{j=0}^{n_0} \left(\frac{1}{j+1} \right)^j - \sum_{j=n_0+1}^{n_1+1} \left(\frac{1}{j+1} \right)^j + \dots \right] \begin{pmatrix} 0.6 \\ 0.8 \end{pmatrix} \right\|_2 < \left\| \begin{pmatrix} 4.5 \\ 6 \end{pmatrix} \right\|_2$$

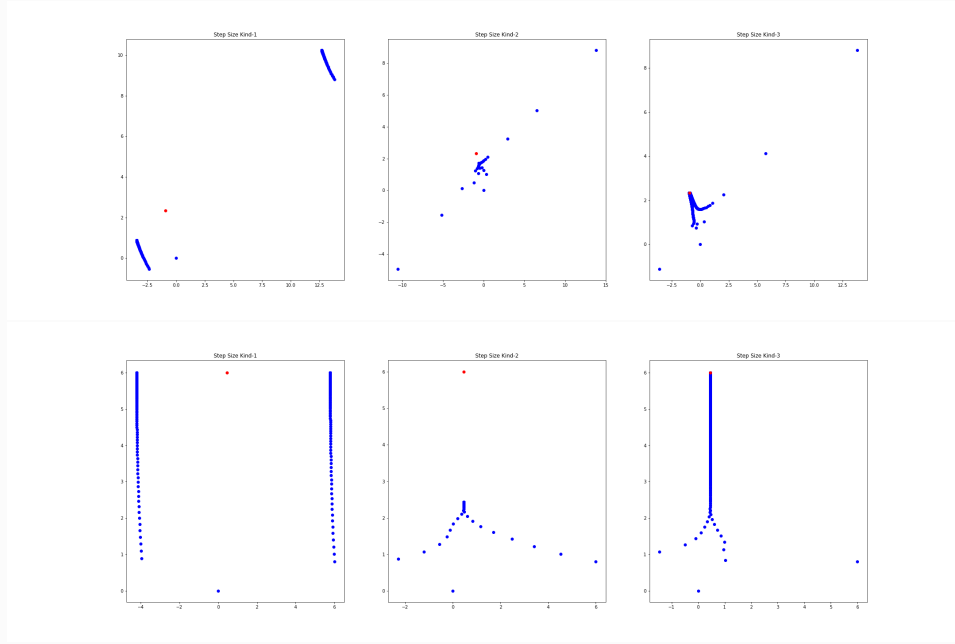
Repeat there preccedures, we can get that the sequence $\{n_0, n_1, \dots\}$

$$\left\| \left[\sum_{j=0}^{n_0} \left(\frac{1}{j+1} \right)^j - \sum_{j=n_0+1}^{n_1} \left(\frac{1}{j+1} \right)^j + \dots \right] \begin{pmatrix} 0.6 \\ 0.8 \end{pmatrix} \right\|_2 \rightarrow \left\| \begin{pmatrix} 4.5 \\ 6 \end{pmatrix} \right\|_2$$

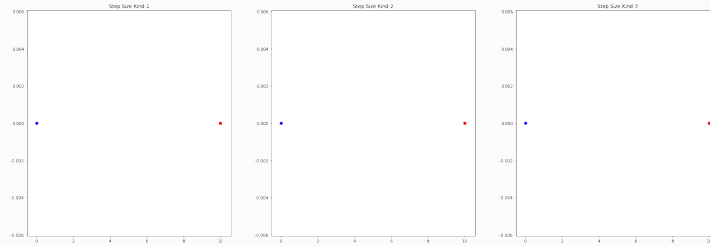
since the remainder of series is infinite. Or in one of n_i , we approach the optimal solution.

(e)

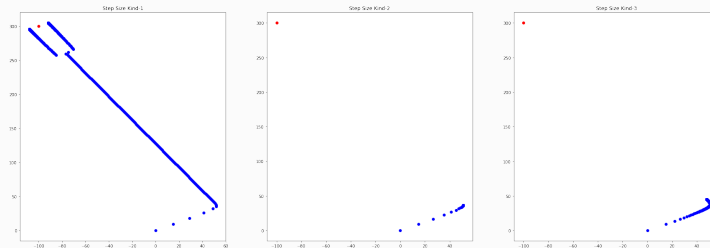
For the given A , only the 3th kind of steps choices finds the optimal solution.



For different choices of b , the results are different. For the first choice of A and $b = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$, all kinds of steps choices work.



For the second choice of A and $b = \begin{pmatrix} 900 \\ 900 \end{pmatrix}$, all kinds of steps choices don't work.



Question 3

(a)

\therefore

$$\begin{aligned}\nabla f(x) &= \frac{1}{2} \nabla \|Ax - b\|_2^2 \\ &= \frac{1}{2} \nabla (Ax - b)^T (Ax - b) \\ &= \frac{1}{2} \nabla [x^T A^T Ax - b^T Ax - x^T A^T b + b^T b] \\ &= A^T Ax - A^T b\end{aligned}$$

given $b = \vec{0}$,

$$\nabla f(x) = A^T Ax$$

\therefore

$$\begin{aligned}x_{i+1} &= x_i - \gamma \nabla f(x_i) \\ &= x_i - \gamma A^T Ax_i \\ &= (I - \gamma A^T A)x_i\end{aligned}$$

with the initial condition x_0 .

(b)

\therefore

$$x_i = (I - \gamma A^T A)^i x_0$$

We want $\forall x_0$, and $\epsilon_0 = x_0 - x$ where x is the optimal solution we are looking for, the following condition is satisfied,

$$\begin{aligned}\epsilon_k &= x_k - x \\ &= (I - \gamma A^T A)\epsilon_{k-1} \\ &= (I - \gamma A^T A)^k \epsilon_0 \rightarrow 0\end{aligned}$$

\therefore

$$\lim_{k \rightarrow \infty} (I - \gamma A^T A)^k = 0$$

i.e. x_i won't blow up arbitrarily over time.

(c)

\therefore

$$\begin{aligned}\varphi(x) &= x - \gamma \nabla f(x) \\ &= x - \gamma (A^T Ax - A^T b) \\ &= (I - \gamma A^T A)x + \gamma (A^T b)\end{aligned}$$

Solution (cont.)

and is a $d \times d$ symmetric matrix, where $a_i \in \mathbb{R}^d$

\therefore

$$I - \gamma A^T A = P^{-1} J P$$

where P is an orthonormal matrix and $J = \begin{pmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & \lambda_d \end{pmatrix}$ is a diagonal matrix

\therefore

$$\begin{aligned} \|\varphi(x) - \varphi(x')\|_2 &= \|(I - \gamma A^T A)(x - x')\|_2 \\ &= \|P^T J P(x - x')\|_2 \\ &= \sqrt{(x - x')^T P^T J P P^T J P (x - x')} \\ &= \sqrt{(x - x')^T P^T J^2 P (x - x')} \\ &\stackrel{y=P(x-x')}{=} \sqrt{y^T J^2 y} \\ &= \sqrt{\sum_{i=1}^n \lambda_i^2 y_i^2} \\ &\leq \max_i \{|\lambda_i|\} \sqrt{\sum_{i=1}^n y_i^2} \\ &= |\lambda_{\max}(I - \gamma A^T A)| \|y\|_2 \end{aligned}$$

\therefore

$$\lambda_{\min}(I - \gamma A^T A) \leq \lambda_{\max}(I - \gamma A^T A) \leq |\lambda_{\max}(I - \gamma A^T A)|$$

\therefore

$$|\lambda_{\max}(I - \gamma A^T A)| \leq \max\{|\lambda_{\min}(I - \gamma A^T A)|, |\lambda_{\max}(I - \gamma A^T A)|\}$$

\therefore

$$\|\varphi(x) - \varphi(x')\|_2 \leq \beta \|x - x'\|_2$$

(d)

\therefore

$$x^* = \arg \min_{x \in \mathbb{R}^d} f(x)$$

\therefore

$$\nabla f(x^*) = 0$$

\therefore

$$\begin{aligned} \|\varphi(x_k) - \varphi(x^*)\|_2 &= \|x_k - \gamma(A^T A x_k - A^T b) - x^*\|_2 \\ &= \|(I - \gamma A^T A)x_k - x^*\|_2 \\ &= \|x_{k+1} - x^*\|_2 \end{aligned}$$

Solution (cont.)

\therefore

$$\begin{aligned}\|x_{k+1} - x^*\|_2 &= \|\varphi(x_k) - \varphi(x^*)\|_2 \\ &\leq \beta \|x_k - x^*\|_2\end{aligned}$$

\therefore

$$\|x_{k+1} - x^*\|_2 \leq \beta^{k+1} \|x_0 - x^*\|_2$$

(e)

\therefore x^* is the minimizer and

$$\nabla f(x) = A^T Ax - A^T b$$

\therefore x^* equals to the solution to $\nabla f(x) = 0$, i.e. $A^T Ax^* = A^T b$

i.e.

$$x^* = (A^T A)^{-1} A^T b$$

\therefore

$$\begin{aligned}f(x^*) &= \frac{1}{2} \|Ax^* - b\|_2^2 \\ &= \frac{1}{2} (Ax^* - b)^T (Ax^* - b) \\ &= \frac{1}{2} [x^{*T} A^T Ax^* - b^T Ax^* - x^{*T} A^T b + b^T b] \\ &= \frac{1}{2} [b^T A (A^T A)^{-1} (A^T A) (A^T A)^{-1} A^T b - 2b^T A (A^T A)^{-1} A^T b + b^T b] \\ &= \frac{1}{2} [-b^T A (A^T A)^{-1} A^T b + b^T b] \\ &= \frac{1}{2} [-b^T A (A^T A)^{-1} (A^T A) (A^T A)^{-1} A^T b + b^T b] \\ &= \frac{1}{2} (-x^{*T} A^T Ax^* + b^T b)\end{aligned}$$

\therefore

$$\begin{aligned}f(x) &= \frac{1}{2} \|Ax - b\|_2^2 \\ &= \frac{1}{2} [(Ax - b)^T (Ax - b)] \\ &= \frac{1}{2} [x^T A^T Ax - b^T Ax - x^T A^T b + b^T b] \\ &= \frac{1}{2} [(x - x^*)^T A^T A (x - x^*) - b^T Ax - x^T A^T b + x^T A^T Ax^* + x^{*T} A^T Ax - x^{*T} A^T Ax^* + b^T b] \\ &= \frac{1}{2} [\|A(x - x^*)\|_2^2 - b^T Ax - x^T A^T b + x^T A^T b + bA^T Ax - x^{*T} A^T Ax^* + b^T b] \\ &= \frac{1}{2} [\|A(x - x^*)\|_2^2 - x^{*T} A^T Ax^*] \\ &= \frac{1}{2} \|A(x - x^*)\|_2^2 - f(x^*)\end{aligned}$$

\therefore

$$f(x) - f(x^*) = \frac{1}{2} \|A(x - x^*)\|_2^2$$

(f)

$\therefore A^T A$ is a symmetric matrix

$$\therefore \exists P, J \in \mathbb{R}^{d \times d} \text{ where } P \text{ is orthonormal, } P^T = P^{-1} \text{ and } J = \begin{pmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_d \end{pmatrix}$$

\therefore

$$\begin{aligned} \|A(x - x^*)\|_2^2 &= (x - x^*)^T A^T A (x - x^*) \\ &= (x - x^*)^T P^T J P (x - x^*) \\ &\stackrel{y=P(x-x^*)}{=} y^T J y \\ &= \sum_{i=1}^d \lambda_i y_i^2 \\ &\leq \max_i \{\lambda_i\} \|y\|_2^2 \\ &= \max_i \{\lambda_i\} \|P(x - x^*)\|_2^2 \\ &= \lambda_{\max}(A^T A) \|x - x^*\|_2^2 \\ &= \alpha \|x - x^*\|_2^2 \end{aligned}$$

\therefore

$$f(x_k) - f(x^*) \leq \frac{\alpha}{2} \|x_k - x^*\|_2^2$$

From (e) we have

$$\begin{aligned} f(x_k) - f(x^*) &\leq \frac{\alpha}{2} (\beta^k \|x_k - x^*\|_2)^2 \\ &= \frac{\alpha}{2} \beta^{2k} \|x_k - x^*\|_2^2 \end{aligned}$$

(g)

$\therefore A^T A$ is positive definite

\therefore

$$\lambda_{\max}(A^T A) \geq \lambda_{\min}(A^T A) > 0$$

\therefore

$$\beta = \max\{|1 - \gamma \lambda_{\max}(A^T A)|, |1 - \gamma \lambda_{\min}(A^T A)|\}$$

Solution (cont.)

$\therefore \forall \epsilon > 0, \exists n \in \mathbb{N}$, s.t. $0 < \left(\frac{\lambda_{\min}(A^T A)}{\lambda_{\max}(A^T A)} \right)^n < \epsilon$, let $\gamma = \frac{1 - \left(\frac{\lambda_{\min}(A^T A)}{\lambda_{\max}(A^T A)} \right)^n}{\lambda_{\min}(A^T A)}$, we have

$$\begin{aligned}
1 - \gamma \lambda_{\min}(A^T A) &= \left(\frac{\lambda_{\min}(A^T A)}{\lambda_{\max}(A^T A)} \right)^n < \epsilon \\
1 - \gamma \lambda_{\min}(A^T A) &\geq 1 - \gamma \lambda_{\max}(A^T A) \\
&= 1 - \left[1 - \left(\frac{\lambda_{\min}(A^T A)}{\lambda_{\max}(A^T A)} \right)^n \right] \frac{\lambda_{\max}(A^T A)}{\lambda_{\min}(A^T A)} \\
&\geq 1 - \left(1 + \frac{\lambda_{\min}(A^T A)}{\lambda_{\max}(A^T A)} \right) \frac{\lambda_{\max}(A^T A)}{\lambda_{\min}(A^T A)} \\
&> -\epsilon \\
1 - \gamma \lambda_{\max}(A^T A) &= 1 - \left[1 - \left(\frac{\lambda_{\min}(A^T A)}{\lambda_{\max}(A^T A)} \right)^n \right] \frac{\lambda_{\max}(A^T A)}{\lambda_{\min}(A^T A)} \\
&\leq 1 - \left[1 - \left(\frac{\lambda_{\min}(A^T A)}{\lambda_{\max}(A^T A)} \right)^n \right] \\
&= \left(\frac{\lambda_{\min}(A^T A)}{\lambda_{\max}(A^T A)} \right)^n < \epsilon
\end{aligned}$$

i.e.

$$\beta < \epsilon$$

When picking $\gamma = \frac{1 - \left(\frac{\lambda_{\min}(A^T A)}{\lambda_{\max}(A^T A)} \right)^n}{\lambda_{\min}(A^T A)}$, we have

$$\begin{aligned}
f(x) - f(x^*) &\leq \frac{\alpha}{2} \beta^{2k} \|x_0 - x^*\|_2^2 \\
&= \frac{\alpha}{2} \left(\frac{\lambda_{\min}(A^T A)}{\lambda_{\max}(A^T A)} \right)^{2nk} \|x_0 - x^*\|_2^2 \\
&= \frac{\alpha}{2} Q^{2nk} \|x_0 - x^*\|_2^2
\end{aligned}$$

Question 4

(a)

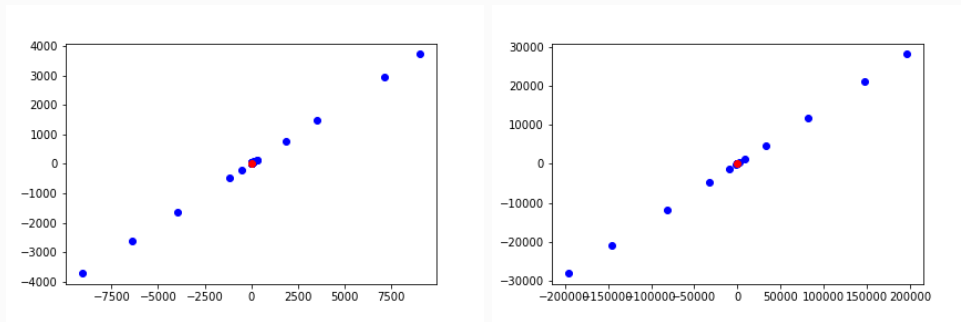
$$L(x_1, y_1) = - \sum_{i=1}^7 (\sqrt{(a_i - x_1)^2 + (b_i - y_1)^2} - d_i)^2$$

From the chain rule,

$$\frac{\partial L}{\partial x_1} = 2 \sum_{i=1}^7 (\sqrt{(a_i - x_1)^2 + (b_i - y_1)^2} - d_i) \frac{(a_i - x_1)}{\sqrt{(a_i - x_1)^2 + (b_i - y_1)^2}}$$

$$\frac{\partial L}{\partial x_2} = 2 \sum_{i=1}^7 (\sqrt{(a_i - x_1)^2 + (b_i - y_1)^2} - d_i) \frac{(b_i - y_1)}{\sqrt{(a_i - x_1)^2 + (b_i - y_1)^2}}$$

(b)



The optimal solutions of these two approaches are

$$\begin{pmatrix} 43.07188566 & 32.71217991 \end{pmatrix} \quad \begin{pmatrix} 43.0718848 & 32.7121788 \end{pmatrix}$$

Starting at a random point seems not as stable as the the former. Sometimes it takes less time to find the optimal solution while sometimes it will take longer time. I choose steps $\gamma_i = \frac{1}{i+1}$

```

1 def gradient(obj_loc , d, sen_loc):
2     g = 2*np.dot(((np.linalg.norm(obj_loc-sen_loc , axis=1)-d
3         )/np.linalg.norm(obj_loc-sen_loc , axis=1)).T,(obj_loc-sen_loc))
4     if np.any(np.isnan(g)):
5         return 0
6     else:
7         return g
8
9 def gradient_descend_step(obj_loc , d, sen_loc , step_count ,
10     step_size):
11     obj_loc = obj_loc - step_size(step_count) * gradient(obj_loc ,
12         d, sen_loc)
13     return np.array(obj_loc)

```

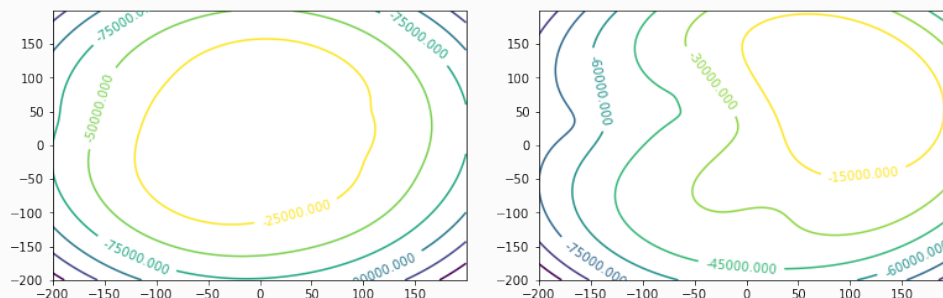
Solution (cont.)

```
14
15 def gradient_descend(obj_loc , d, sen_loc , total_step_count ,
16     step_size , err=0):
17     positions = [np.array(obj_loc)]
18     for k in range(total_step_count):
19         positions.append(gradient_descend_step(positions[-1],
20             d, sen_loc , k, step_size))
21
22     return np.array(positions)
23
24 initial_position = np.array([0, 0]) # position at iteration 0
25 total_step_count = 1000 # number of GD steps to take
26 step_size = lambda i:1/(i+1) # step size at iteration i
27 err = 1e-15
28 positions = gradient_descend(initial_position , single_distance ,
29     sensor_loc , total_step_count , step_size , err)
30 print(positions)
31 plt.scatter(positions[:,0], positions[:,1], c='blue')
32 plt.scatter(obj_loc[:,0], obj_loc[:,1], c='red')
33 plt.plot()
34 plt.savefig('4b1.png')
35 plt.show()
36
37 initial_position = 100*np.random.randn(2)
38 print(initial_position)
39 total_step_count = 1000 # number of GD steps to take
40 step_size = lambda i: 1/(i+1) # step size at iteration i
41 err = 1e-15
42 positions = gradient_descend(initial_position , single_distance ,
43     sensor_loc , total_step_count , step_size , err)
44 print(positions)
45 plt.scatter(positions[:,0], positions[:,1], c='blue')
46 plt.scatter(obj_loc[:,0], obj_loc[:,1], c='red')
47 plt.plot()
48 plt.savefig('4b2.png')
49 plt.show()
```

(c)

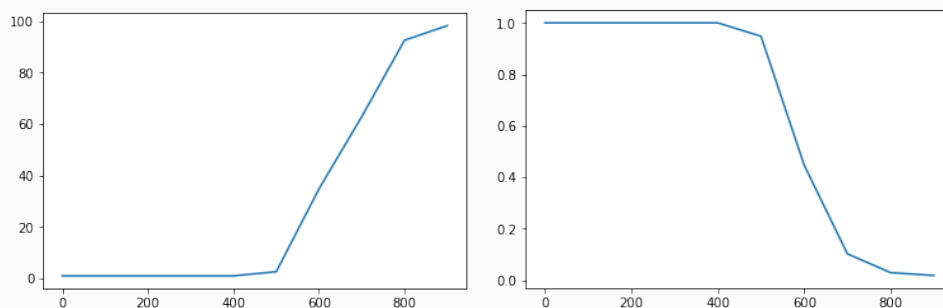
Contour plot for $x_1 = (0,0)$ and $x_1 = (100,100)$:

Solution (cont.)



The average number of local minima against x_1 :

The average proportion against x_1 :



When the $\|x_i\|_2$ is big enough, the easier to find the local minima. However, it will find more local minima than the global minima.

```
1 def gradient(obj_loc , d, sen_loc):
2     g = 2*np.dot(((np.linalg.norm(obj_loc-sen_loc , axis=1)-d)
3         /np.linalg.norm(obj_loc-sen_loc , axis=1)).T,( obj_loc-sen_loc))
4     if np.any(np.isnan(g)):
5         return 0
6     else:
7         return g
8
9 def gradient_descend_step(obj_loc , d, sen_loc , step_count ,
10     step_size):
11     obj_loc = obj_loc - step_size(step_count) * gradient(obj_loc , d,
12         sen_loc)
13     return np.array(obj_loc)
14
15 def gradient_descend(obj_loc , d, sen_loc , total_step_count ,
16     step_size , err=0):
17     positions = np.array(obj_loc)
18     for k in range(total_step_count):
19         positions = gradient_descend_step(positions , d,
20             sen_loc , k, step_size)
21
```

Solution (cont.)

```
22     return positions
23
24 def generate_data_given_location(sensor_loc , obj_loc , k = 7,
25     d = 2, sigma=1):
26     assert k, d == sensor_loc.shape
27
28     distance = scipy.spatial.distance.cdist(obj_loc ,
29         sensor_loc ,
30         metric='euclidean')
31     distance += np.random.randn(1, k)*sigma
32     return distance
33
34 def log_likelihood(obj_loc , sensor_loc , distance):
35     diff_distance = np.sqrt(np.sum((sensor_loc - obj_loc)**2,
36         axis = 1))- distance
37     func_value = -sum((diff_distance)**2)/2
38     return func_value
39
40 np.random.seed(100)
41 # Sensor locations.
42 sensor_loc = generate_sensors()
43 num_gd_replicates = 100
44
45 # Object locations.
46 obj_locs = [[[i,i]] for i in np.arange(0,1000,100)]
47
48 distances = []
49 for i in range(len(obj_locs)):
50     distances.append(generate_data_given_location(sensor_loc ,
51         obj_locs[i])[0])
52 distances = np.array(distances)
53
54 x_axis = np.arange(-200.0, 200.0, 1)
55 y_axis = np.arange(-200.0, 200.0, 1)
56 X, Y = np.meshgrid(x_axis , y_axis)
57 m ,n = np.shape(X)
58 like = [[log_likelihood([X[j][i],Y[j][i]], sensor_loc ,
59     distances[0][0]) for i in range(n)] for j in range(m)]
60 CS = plt.contour(X, Y, like)
61 plt.clabel(CS, inline=1, fontsize=10)
62 plt.savefig('4c1.png')
63 plt.show()
64
```

Solution (cont.)

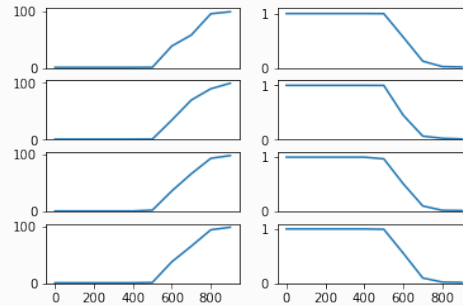
```
65 x_axis = np.arange(-200.0, 200.0, 1)
66 y_axis = np.arange(-200.0, 200.0, 1)
67 X, Y = np.meshgrid(x_axis, y_axis)
68 m, n = np.shape(X)
69 like = [[log_likelihood([X[j][i], Y[j][i]], sensor_loc, distances[1][0])]
70         for i in range(n)] for j in range(m)]
71 CS = plt.contour(X, Y, like)
72 plt.xlabel(CS, inline=1, fontsize=10)
73 plt.savefig('4c2.png')
74 plt.show()
75
76 total_step_count = 1000 # number of GD steps to take
77 step_size = lambda i: 0.1 # step size at iteration i
78 err = 0
79 positions = []
80 for i in range(len(obj_locs)):
81     position = []
82     for j in range(10):
83         optimal_solution = np.zeros((num_gd_replicates, 2))
84         initial_position = np.random.randn(num_gd_replicates, 2) * (
85             obj_locs[i][0][0] + 1)
86         for k in range(num_gd_replicates):
87             optimal_solution[k, :] = gradient_descend(
88                 initial_position[k], distances[i], sensor_loc,
89                 total_step_count, step_size, err)[-1]
90         position += [optimal_solution]
91     print(i)
92     positions += [np.array(position)]
93 positions = np.array(positions)
94
95 result = np.zeros((10, 10))
96 for i in range(10):
97     for j in range(10):
98         result[i, j] = len(np.unique(np.round(positions[i][j], 2),
99                                     axis=0))
100 average_result = np.mean(result, axis=1)
101 plt.plot(np.arange(0, 1000, 100), average_result)
102
103 result2 = np.zeros((10, 10))
104 for i in range(10):
105     for j in range(10):
106         result2[i, j] = max(np.unique(np.round(positions[i][j], 2),
107                                     axis=0, return_counts=True)[1])
```


Solution (cont.)

```
108 average_result2 = np.mean(result2 , axis=1)/ num_gd_replicates
109 plt.plot(np.arange(0,1000,100), average_result2)
```

(d)

We can see that the average number of local minima against x_1 and the average proportion against x_1 remain almost the same when change the variance of the measurement noise.



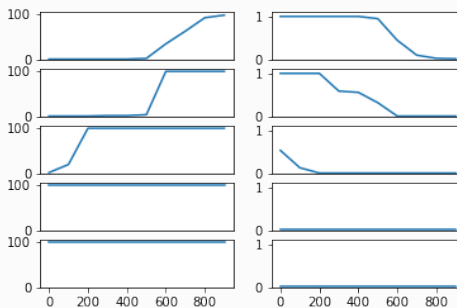
```
1  np.random.seed(100)
2  sensor_loc = generate_sensors()
3  num_gd_replicates = 100
4  sigma = np.arange(1,5)
5
6  obj_locs = [[[i,i]] for i in np.arange(0,1000,100)]
7
8  total_step_count = 1000 # number of GD steps to take
9  step_size = lambda i: 0.1 # step size at iteration i
10 err = 0
11 positions2 = []
12
13 for e in range(len(sigma)):
14     distance2 = []
15     for i in range(len(obj_locs)):
16         distance2.append(generate_data_given_location(sensor_loc ,
17             obj_locs[i] , sigma=10^(-i))[0])
18     distance2 = np.array(distance2)
19
20     position2 = []
21     for i in range(len(obj_locs)):
22         p = []
23         for j in range(10):
24             optimal_solution2 = np.zeros((num_gd_replicates,2))
25             initial_position2 = np.random.randn(num_gd_replicates,2)
26             *(obj_locs[i][0][0]+1)
```

Solution (cont.)

```
27         for k in range(num_gd_replicates):
28             optimal_solution2[k] = gradient_descend(
29                 initial_position2[k], distance2[i], sensor_loc,
30                 total_step_count, step_size, err)[-1]
31             p += [optimal_solution2]
32             position2 += [np.array(p)]
33         position2 = np.array(position2)
34         positions2.append(position2)
35
36     plt.subplot(len(sigma), 2, 2*e+1)
37     result21 = np.zeros((10,10))
38     for i in range(10):
39         for j in range(10):
40             result21[i,j] = len(np.unique(np.round(position2[i][j],
41                 2), axis=0))
42     average_result21 = np.mean(result21, axis=1)
43     plt.plot(np.arange(0,1000,100), average_result21)
44     plt.subplot(len(sigma), 2, 2*e+2)
45     result22 = np.zeros((10,10))
46     for i in range(10):
47         for j in range(10):
48             result22[i,j] = max(np.unique(np.round(position2[i][j],
49                 2), axis=0, return_counts=True)[1])
50     average_result22 = np.mean(result22, axis=1)/ num_gd_replicates
51     plt.plot(np.arange(0,1000,100), average_result22)
52     print(e)
53     plt.show()
```

(e)

We can see that the average number of local minima against x_1 increases and the average proportion against x_1 decreases when number of sensor increases. I.e, when increasing the number of sensors, it is much easier to find the local minima but harder to find the global minima.



Solution (cont.)

```
1  np.random.seed(100)
2  # Sensor locations.
3  num_gd_replicates = 100
4  num_sensor = np.arange(7,28,5)
5
6  # Object locations.
7  obj_locs = [[[i,i]] for i in np.arange(0,1000,100)]
8
9  total_step_count = 1000 # number of GD steps to take
10 step_size = lambda i: 0.1 # step size at iteration i
11 err = 0
12 positions3 = []
13
14 for e in range(len(num_sensor)):
15     sensor_loc = generate_sensors(k=num_sensor[e])
16     distance3 = []
17     for i in range(len(obj_locs)):
18         distance3.append(generate_data_given_location(sensor_loc ,
19                                                         obj_locs[i], k=num_sensor[e])[0])
20     distance3 = np.array(distance3)
21
22     position3 = []
23     for i in range(len(obj_locs)):
24         p = []
25         for j in range(10):
26             optimal_solution3 = np.zeros((num_gd_replicates,2))
27             initial_position3 = np.random.randn(num_gd_replicates,2)
28                 *(obj_locs[i][0][0]+1)
29             for k in range(num_gd_replicates):
30                 optimal_solution3[k] = gradient_descend(
31                     initial_position3[k], distance3[i], sensor_loc ,
32                     total_step_count , step_size , err)
33             p += [optimal_solution3]
34         position3 += [np.array(p)]
35     position3 = np.array(position3)
36     positions3.append(position3)
37
38     plt.subplot(len(num_sensor),2,2*e+1)
39     result31 = np.zeros((10,10))
40     for i in range(10):
41         for j in range(10):
42             result31[i,j] = len(np.unique(np.round(position3[i][j],2
43                                                     ),axis=0))
```

Solution (cont.)

```

44     average_result31 = np.mean(result31 , axis=1)
45     plt.plot(np.arange(0,1000,100), average_result31)
46     plt.subplot(len(num_sensor),2,2*e+2)
47     result32 = np.zeros((10,10))
48     for i in range(10):
49         for j in range(10):
50             result32[i,j] = max(np.unique(np.round(position3[i][j],2
51                                     ), axis=0, return_counts=True)[1])
52     average_result32 = np.mean(result32 , axis=1)/ num_gd_replicates
53     plt.plot(np.arange(0,1000,100), average_result32)
54     print(e)
55 np.save('4e', positions3)
56 plt.show()

```

(f)

$$L(a_1, \dots, a_6, b_1, \dots, b_7) = - \sum_{i=1}^7 \sum_{j=1}^{100} (\sqrt{(a_i - x_j)^2 + (b_i - y_j)^2} - d_{ij})^2$$

From the chain rule,

$$\frac{\partial L}{\partial a_i} = -2 \sum_{j=1}^{100} (\sqrt{(a_i - x_j)^2 + (b_i - y_j)^2} - d_{ij}) \frac{(a_i - x_j)}{\sqrt{(a_i - x_j)^2 + (b_i - y_j)^2}}$$

$$\frac{\partial L}{\partial b_i} = -2 \sum_{j=1}^{100} (\sqrt{(a_i - x_j)^2 + (b_i - y_j)^2} - d_{ij}) \frac{(b_i - y_j)}{\sqrt{(a_i - x_j)^2 + (b_i - y_j)^2}}$$

DataSet	Train	Test1	Test2
MSE	1469.01800279	10163.56297	657407.238113

```

1  np.random.seed(100)
2  # Sensor locations.
3  sensor_loc = generate_sensors()
4
5  train = np.random.randn(100,2)*100
6  test1 = np.random.randn(100,2)*100
7  test2 = 300+np.random.randn(100,2)*100
8
9  train_distances = []
10 for i in range(len(train)):
11     train_distances.append(generate_data_given_location(sensor_loc ,
12                 train[i,np.newaxis])[0])
13 train_distances = np.array(train_distances)

```

Solution (cont.)

```
14 test1_distances = []
15 for i in range(len(test1)):
16     test1_distances.append(generate_data_given_location(sensor_loc ,
17         test1[i,np.newaxis])[0])
18 test1_distances = np.array(test1_distances)
19 test2_distances = []
20 for i in range(len(test2)):
21     test2_distances.append(generate_data_given_location(sensor_loc ,
22         test2[i,np.newaxis])[0])
23 test2_distances = np.array(test2_distances)
24
25 def gradient1(obj_loc , d, sen_loc):
26     g = np.zeros_like(sen_loc)
27     for i in range(len(g)):
28         g[i] = -2*np.dot(((np.linalg.norm(obj_loc-sen_loc[i], axis=1)
29             -d[:,i])/np.linalg.norm(obj_loc-sen_loc[i], axis=1)).T,
30             (obj_loc-sen_loc[i]))
31     if np.any(np.isnan(g)):
32         return 0
33     else:
34         return g
35
36 def gradient_descend_step1(obj_loc , d, sen_loc , step_count ,
37     step_size):
38     sen_loc = sen_loc - step_size(step_count) * gradient1(obj_loc ,
39         d, sen_loc)
40     return sen_loc
41
42 def gradient_descend1(obj_loc , d, sen_loc , total_step_count ,
43     step_size , err=0):
44     positions = [np.array(sen_loc)]
45     for k in range(total_step_count):
46         new = gradient_descend_step1(obj_loc , d, positions[-1],
47             k, step_size)
48         if np.linalg.norm(positions[-1]-new)<err:
49             break
50         else:
51             positions.append(new)
52     return np.array(positions)
53
54 position4 = []
55 total_step_count = 10000 # number of GD steps to take
56 step_size = lambda i: 1/(1+i)#0.001 # step size at iteration i
```

Solution (cont.)

```
57 err = 1e-10
58 for i in range(num_gd_replicates):
59     p = []
60     initial_position4 = np.random.randn(7,2)*100
61     optimal_solution4 = gradient_descend1(train, train_distances,
62         initial_position4, total_step_count, step_size, err)
63     p += [optimal_solution4[-1]]
64     position4 += [np.array(p)]
65 position4 = np.array(position4)
66 result4 = np.unique(np.round(position4,2), axis=0)[np.argmax(
67     np.unique(np.round(position4,2), axis=0, return_counts=True)[1])]
68
69 def gradient2(obj_loc, d, sen_loc):
70     g = 2*np.dot((((np.linalg.norm(obj_loc-sen_loc, axis=1)-d)/
71         np.linalg.norm(obj_loc-sen_loc, axis=1))).T, (obj_loc-sen_loc))
72     if np.any(np.isnan(g)):
73         return 0
74     else:
75         return g/len(g)
76
77 def gradient_descend_step2(obj_loc, d, sen_loc, step_count,
78     step_size):
79     obj_loc = obj_loc - step_size(step_count) * gradient2(obj_loc,
80         d, sen_loc)
81     return obj_loc
82
83 def gradient_descend2(obj_loc, d, sen_loc, total_step_count,
84     step_size, err=0):
85     positions = [np.array(obj_loc)]
86     for k in range(total_step_count):
87         new = gradient_descend_step2(positions[-1], d, sen_loc,
88             k, step_size)
89         if np.max(np.linalg.norm(positions[-1]-new, axis=1))<err:
90             break
91         else:
92             positions.append(new)
93     return positions
94
95 total_step_count = 1000 # number of GD steps to take
96 step_size = lambda i: 1/(1+i)#0.001 # step size at iteration i
97 err = 1e-10
98
99 result5 = np.zeros_like(train)
```

Solution (cont.)

```
100 for j in range(len(train)):
101     p = []
102     for i in range(num_gd_replicates):
103         initial_position5 = np.random.randn(1,2)*100
104         optimal_solution5 = gradient_descend2(initial_position5 ,
105             train_distances[j], result4[0], total_step_count ,
106             step_size , err)
107         p += [optimal_solution5[-1]]
108     position5 = np.array(p)
109     result5[j] = np.unique(np.round(position5,2), axis=0)[np.argmax
110         (np.unique(np.round(position5,2), axis=0, return_counts=True)
111             [1])]
112
113 result51 = np.zeros_like(test1)
114 for j in range(len(test1)):
115     p = []
116     for i in range(num_gd_replicates):
117         initial_position51 = np.random.randn(1,2)*100
118         optimal_solution51 = gradient_descend2(initial_position51 ,
119             test1_distances[j], result4[0], total_step_count ,
120             step_size , err)
121         p += [optimal_solution51[-1]]
122     position51 = np.array(p)
123     result51[j] = np.unique(np.round(position51,2), axis=0)[
124         np.argmax(np.unique(np.round(position51,2), axis=0,
125             return_counts=True)[1])]
126
127 result52 = np.zeros_like(test2)
128 for j in range(len(test2)):
129     p = []
130     for i in range(num_gd_replicates):
131         initial_position52 = np.random.randn(1,2)*100
132         optimal_solution52 = gradient_descend2(initial_position52 ,
133             test2_distances[j], result4[0], total_step_count ,
134             step_size , err)
135         p += [optimal_solution52[-1]]
136     position52 = np.array(p)
137     result52[j] = np.unique(np.round(position52,2), axis=0)[
138         np.argmax(np.unique(np.round(position52,2), axis=0,
139             return_counts=True)[1])]
140
141 MSE_train = np.sum(np.linalg.norm(result5 - train, axis=1)**2)
142 /len(train)
```

Solution (cont.)

```
143 MSE_test1 = np.sum(np.linalg.norm(result51 - test1 , axis=1)**2)
144     /len(test2)
145 MSE_test2 = np.sum(np.linalg.norm(result52 - test2 , axis=1)**2)
146     /len(test1)
```


Question 5

Have uploaded to Gradescope.

Question 6

Question What is the difference between various kinds of gradient descend method?

Solution

For cost function θ , data set D with m data, η is the learning rate

(1) Batch Gradient Descend (BGD)

$$\theta_j = \theta_j - \eta \frac{1}{m} \sum_{i=1}^m [h_{\theta}(x_i) - y_i] x_i$$

It cannot work well when m is very large.

(2) Mini-Batch Gradient Descend (MBGD)

Given batch size b , we only calculate the gradient for a small subset of D . We can split D into several subset or we can just choose each S randomly.

$$\theta_j = \theta_j - \eta \frac{1}{b} \sum_{\substack{S \subset D \\ |S|=b}} [h_{\theta}(x_i) - y_i] x_i$$

(3) Stochastic Gradient Descend (SGD)

For $i = 1, 2, \dots, m$,

$$\theta_j = \theta_j - \eta [h_{\theta}(x_i) - y_i] x_i$$

it may not return a most accuracy optimal solution but it is useful because it cost less time.

Method	BGD	MBGD	SGD
Accuracy	high	medium	low
Time consuming	low	medium	high