

In the name of God

## Convex Optimization Project

Dr. Yassaee

Mobin Khatib

99106114

Amir Ali Loghmani

99102145



# 1 Theory Question

Given the function

$$f(x) = \frac{1}{2}x^T Ax - b^T x + c \quad (1)$$

, we use the gradient descent update rule at step : t with learning rate  $\eta$

$$x_{t+1} = x_t - \eta_t \nabla f(x_t) \quad (2)$$

We want to show that:

$$x_{t+1} = (I - \eta_t A)x_t + \eta_t b \quad (3)$$

Let's assume  $f(x)$  has an optimal solution  $x^*$ :

$$x_{t+1} = (I - \eta_t A)x_t + \eta_t b \quad (4)$$

$$x_{t+1} - x^* = (x_t - \eta_t(Ax_t - b)) - (x^* - \eta_t(Ax^* - b)) \quad (5)$$

$$x_{t+1} - x^* = (I - \eta_t A)(x_t - x^*) = \prod_{k=1}^n (I - \eta_k A)(x_1 - x^*) \quad (6)$$

Considering above equation The convergence of gradient descent for quadratic functions can be analyzed mathematically. In particular, for a strictly convex quadratic function, the sequence of iterates generated by gradient descent will converge to the optimal solution  $x^*$  as the number of iterations increases. The rate of convergence depends on the eigenvalues of the matrix A and the choice of the learning rate

$$\lim_{T \rightarrow \infty} \|x(T) - x^*\| = 0 \quad (7)$$

When we want this because all of the  $x$ 's from 1 to n are multiplied in equation 6 if one of the answers for  $x_t - x^* = 0$  then equation number 6 becomes zero and equation number 7 will be satisfied

## 2 Theory Question

$$x_{(t+1)} = \arg \min_{x \in R^n} f(x_{(t)}) + \langle \nabla f(x_{(t)}), x - x_{(t)} \rangle + \frac{1}{2\eta} \|x - x_{(t)}\|^2 \quad (8)$$

where:

$f(x_{(t)})$  represents the objective function evaluated at  $x_{(t)}$ .  $\nabla f(x_{(t)})$  denotes the gradient of the objective function with respect to  $x$ , evaluated at  $x_{(t)}$ .  $\langle \nabla f(x_{(t)}), x - x_{(t)} \rangle$  represents the inner product (dot product) between the gradient  $\nabla f(x_{(t)})$  and the difference  $x - x_{(t)}$ .  $\frac{1}{2\eta} \|x - x_{(t)}\|^2$  is a regularization term that penalizes the distance between  $x$  and  $x_{(t)}$ , scaled by the parameter  $\eta$ . To find  $x_{(t+1)}$ , we need to compute the value of  $x$  that minimizes this expression. This can be done by taking the derivative of the expression with respect to  $x$ , setting it equal to zero, and solving for  $x$ . Let's differentiate the expression:

$$\frac{d}{dx} f(x_{(t)}) + \langle \nabla f(x_{(t)}), x - x_{(t)} \rangle + \frac{1}{2\eta} \|x - x_{(t)}\|^2 = 0 \quad (9)$$

To simplify the calculation, let's introduce some intermediate variables:

Let  $A = \nabla f(x_{(t)})$ , which is the gradient of the objective function. Let  $B = x - x_{(t)}$ , which is the difference between  $x$  and  $x_{(t)}$ . Now, let's differentiate each term of the expression:

since  $f(x_{(t)})$  is a constant with respect to  $x$  then:

$$\frac{d}{dx} f(x_{(t)}) = 0 \quad (10)$$

$$\frac{d}{dx} \langle \nabla f(x_{(t)}), x - x_{(t)} \rangle = \langle \nabla f(x_{(t)}), \frac{d}{dx} (x - x_{(t)}) \rangle = \quad (11)$$

$$\langle \nabla f(x_{(t)}), dx \rangle = \langle A, dx \rangle. \quad (12)$$

$$\frac{d}{dx} \frac{1}{2\eta} \|x - x_{(t)}\|^2 = \frac{1}{\eta} \frac{d}{dx} \|B\|^2 = \quad (13)$$

$$\frac{1}{\eta} \frac{d}{dx} \langle B, B \rangle = \frac{1}{\eta} \cdot 2 \langle B, \frac{dB}{dx} \rangle = \quad (14)$$

$$\frac{2}{\eta} \langle B, \frac{dB}{dx} \rangle. \quad (15)$$

Putting it all together:

$$\bullet = \frac{d}{dx}f(x_{(t)}) + \langle \nabla f(x_{(t)}), x - x_{(t)} \rangle + \frac{1}{2\eta} \|x - x_{(t)}\|^2 = \quad (16)$$

$$\langle A, dx \rangle + \frac{1}{\eta} \langle B, \frac{dB}{dx} \rangle = \quad (17)$$

$$\langle A, dx \rangle + \frac{1}{\eta} \langle B, (dx - dx_{(t)}) \rangle = \quad (18)$$

$$\langle A, dx \rangle + \frac{1}{\eta} \langle B, dx \rangle = \quad (19)$$

$$\langle A + \frac{1}{\eta} B, dx \rangle. \quad (20)$$

To satisfy the equation for all values of  $dx$ , the coefficient of  $dx$  must be zero:

$$A + \frac{1}{\eta} B = \bullet. \quad (21)$$

Substituting the values of  $A$  and  $B$ , we get:

$$\nabla f(x_{(t)}) + \frac{1}{\eta} (x - x_{(t)}) = \bullet. \quad (22)$$

Rearranging the equation, we find:

$$\nabla f(x_{(t)}) = -\frac{1}{\eta} (x - x_{(t)}). \quad (23)$$

Now, solving for  $x$ :

$$x = x_{(t)} - \eta \nabla f(x_{(t)}). \quad (24)$$

Therefore, the value of  $x$  that minimizes the given expression is given by:

$$x_{(t+1)} = x_{(t)} - \eta \nabla f(x_{(t)}). \quad (25)$$

This expression represents an update rule commonly used in optimization algorithms like gradient descent, where  $x_{(t+1)}$  is updated based on the gradient of the objective function at  $x_{(t)}$ , scaled by a learning rate  $\eta$ .

The equation suggests that we can improve our current estimate or solution  $x(t)$  by taking a step in the opposite direction of the gradient  $\nabla f(x(t))$ . The gradient  $\nabla f(x(t))$  provides information about the direction of steepest ascent of the objective function  $f(x)$ , so moving in the opposite direction corresponds to descending or The relationship captures an iterative optimization process minimizing the function. where each step involves updating the current estimate  $x(t)$  based on the gradient information and the regularization term to minimize the objective function  $f(x)$ . The algorithm aims to converge to a solution that optimizes the objective function within the given constraints.

### 3 Theory Question

With replacing  $x_{(t+1)} = x_{(t)} - \eta_t v_t$

$$||x_{t+1} - x^*||_{\Psi}^2 = ||x_t - \eta_t v_t - x^*||_{\Psi}^2 \quad (26)$$

As we have  $\eta_t \geq 0$  and  $v_t \geq 0$ (because it is in  $\sigma f(x_t)$ ) thus:

$$||x_{t+1} - x^*||_{\Psi}^2 \leq ||x_t - x^*||_{\Psi}^2 \quad (27)$$

### 4 Theory Question

We know that we can replace  $x, y$  with any two other things in the L-lipchitz formula and so with this, we should prove that:

$$|f(x) - f(y)| \leq L||x - y|| \leftrightarrow ||v|| \leq L, v \in \sigma f(w) \quad *** \quad (28)$$

Because the left side of inequality \*\*\* is true for any  $x, y$  , therefore we can choose  $x, y$  in the way of that norm of  $z$  become equal or bigger than zero:

$$|x - y| = 1, (x - y)^T z = ||z||_{\Psi}, z \in \sigma f(w) \quad (29)$$

And we know that  $f$  is convex so:

$$f(x) \geq f(y) + (x - y)^T z \quad (30)$$

$$\rightarrow |f(x) - f(y)| \geq z^T(x - y) \quad (31)$$

$$|f(x) - f(y)| \leq L\|x - y\| \quad (32)$$

we have  $\|x - y\| = 1$  and  $(x - y)^T z = \|z\|_2$  so with this above equations:

$$\|z\| \leq L \quad (33)$$

And for right side of the \*\*\* inequality (because  $f$  is convex):

$$f(x) \geq f(y) + (x - y)^T z, z^T \in \partial f(y) \rightarrow |f(x) - f(y)| \geq z^T(x - y) \quad (34)$$

$$\leq \|z\| \|x - y\| \leq L\|x - y\| \rightarrow \quad (35)$$

$$|f(x) - f(y)| \leq L\|x - y\| \quad (36)$$

## 5 Theory Question

$$\|x_{t+1} - x^*\|_{\Psi}^2 = \|x_t - \eta v_t - x^*\|_{\Psi}^2 \quad (37)$$

$$= \|x_t - x^*\|_{\Psi}^2 - 2\eta_t v_t^T(x_t - x^*) + \eta_t^2 \|v_t\|_{\Psi}^2 \quad (38)$$

$$\leq \|x_t - x^*\|_{\Psi}^2 - 2\eta_t(f(x_t) - f^*) + \eta_t^2 \|v_t\|_{\Psi}^2 \quad (39)$$

where  $f^* = f(x^*)$  The last line follows from the definition of subgradient, which gives

$$f(x^*) \geq f(x_t) + v_t^T(x^* - x_t) \quad (40)$$

So we have:

$$\frac{1}{\Psi} \|x_{t+1} - x^*\|_{\Psi}^2 \leq \frac{1}{\Psi} \|x_t - x^*\|_{\Psi}^2 - 2\eta_t(f(x_t) - f(x^*)) + \frac{\eta_t^2}{\Psi} \|v_t\|_{\Psi}^2 \quad (41)$$

## 6 Theory Question

$$\eta_t(f(x_t) - f(x^*)) \leq \frac{1}{\Psi} \|x_t - x^*\|_{\Psi}^2 - \frac{1}{\Psi} \|x_{t+1} - x^*\|_{\Psi}^2 + \frac{1}{\Psi} \eta_t^2 \|v_t\|_{\Psi}^2 \quad (42)$$

Given:

$$\|v\|^{\mathfrak{Y}} \leq \rho \quad (43)$$

and from the result of question 4:

$$\|x^*\|^{\mathfrak{Y}} \leq B \quad (44)$$

By recursively substituting, we have:

$$t = \mathfrak{I} \Rightarrow \|x^* - x_{(\mathfrak{I})}\|^{\mathfrak{Y}} \leq B\eta_{\mathfrak{I}}(f(x_{(\mathfrak{I})}) - f(x^*)) \leq \frac{\mathfrak{I}}{\mathfrak{Y}}B^{\mathfrak{Y}} - \frac{\mathfrak{I}}{\mathfrak{Y}}\|x_{(\mathfrak{I})} - x^*\|^{\mathfrak{Y}} + \frac{\mathfrak{I}}{\mathfrak{Y}}\eta_{\mathfrak{I}}^{\mathfrak{Y}}\rho^{\mathfrak{Y}} \quad (45)$$

$$t = \mathfrak{Y} \Rightarrow \eta_{\mathfrak{Y}}(f(x_{(\mathfrak{Y})}) - f(x^*)) \leq \frac{\mathfrak{I}}{\mathfrak{Y}}B^{\mathfrak{Y}} - \frac{\mathfrak{I}}{\mathfrak{Y}}\|x_{(\mathfrak{Y})} - x^*\|^{\mathfrak{Y}} - \frac{\mathfrak{I}}{\mathfrak{Y}}\|x_{(\mathfrak{Y})} - x^*\|^{\mathfrak{Y}} + \frac{\mathfrak{I}}{\mathfrak{Y}}\eta_{\mathfrak{Y}}^{\mathfrak{Y}}\rho^{\mathfrak{Y}} \quad (46)$$

...

$$t = T \Rightarrow \eta_T(f(x_{(T)}) - f(x^*)) \leq \frac{\mathfrak{I}}{\mathfrak{Y}}B^{\mathfrak{Y}} - \frac{\mathfrak{I}}{\mathfrak{Y}}\|x_{(T)} - x^*\|^{\mathfrak{Y}} - \frac{\mathfrak{I}}{\mathfrak{Y}}\|x_{(T+\mathfrak{I})} - x^*\|^{\mathfrak{Y}} + \frac{\mathfrak{I}}{\mathfrak{Y}}\eta_T^{\mathfrak{Y}}\rho^{\mathfrak{Y}} \quad (47)$$

$$\sum_{t=\mathfrak{I}}^T \eta_t (f(x_{(t)}) - f(x^*)) \leq \frac{\mathfrak{I}}{\mathfrak{Y}}B^{\mathfrak{Y}} + \frac{\mathfrak{I}}{\mathfrak{Y}}\rho^{\mathfrak{Y}} \sum_{t=\mathfrak{I}}^T \eta_t^{\mathfrak{Y}} - \frac{\mathfrak{I}}{\mathfrak{Y}}\|x_{(T+\mathfrak{I})} - x^*\|^{\mathfrak{Y}} \quad (48)$$

Since we have:

$$\frac{\mathfrak{I}}{\mathfrak{Y}}\|x_{(T+\mathfrak{I})} - x^*\|^{\mathfrak{Y}} \geq \cdot \quad (49)$$

From this, we can conclude that:

$$\sum_{t=\mathfrak{I}}^T \eta_t (f(x_{(t)}) - f(x^*)) \leq \frac{\mathfrak{I}}{\mathfrak{Y}}B^{\mathfrak{Y}} + \frac{\mathfrak{I}}{\mathfrak{Y}}\rho^{\mathfrak{Y}} \sum_{t=\mathfrak{I}}^T \eta_t^{\mathfrak{Y}} \quad (50)$$

## 7 Theory Question

We define this:

$$\sum_{t=\mathfrak{I}}^T \eta_t = a \quad (51)$$

$$\sum_{t=1}^T \eta_t f(x_{(t)}) - af(x^*) \leq \frac{1}{\Upsilon} B^\Upsilon + \frac{1}{\Upsilon} \rho^\Upsilon \sum_{t=1}^T \eta_t^\Upsilon \quad (52)$$

Now we divide this by  $a$  and we will have:

$$\sum_{t=1}^T \frac{\eta_t}{a} f(x_{(t)}) - f(x^*) \leq \frac{\frac{1}{\Upsilon} B^\Upsilon}{a} + \frac{\frac{1}{\Upsilon} \rho^\Upsilon \sum_{t=1}^T \eta_t^\Upsilon}{a} \quad (53)$$

On the other hand because  $f$  is convex, we have:

$$f(\bar{x}_T) = f\left(\sum_{t=1}^T \frac{\eta_t}{a} x_{(t)}\right) \leq \sum_{t=1}^T \frac{\eta_t}{a} f(x_{(t)}) \quad (54)$$

From the two above equations, it follows that:

$$f(\bar{x}_T) - f(x^*) \leq \frac{B^\Upsilon + \rho^\Upsilon \sum_{t=1}^T \eta_t^\Upsilon}{\Upsilon \sum_{t=1}^T \eta_t} \quad (55)$$

## 8 Theory Question

with substitution  $\eta_t = \eta$

$$f(\bar{x}_T) - f(x^*) \leq \frac{B^\Upsilon + T\rho^\Upsilon \eta}{\Upsilon T \eta} \quad (56)$$

To obtain the optimal bound, we take the derivative of the right-hand side and set it equal to zero:

$$\frac{d}{d\eta} \left( \frac{B^\Upsilon + T\rho^\Upsilon \eta}{\Upsilon T \eta} \right) = 0 \rightarrow \eta = \frac{B}{\rho} \sqrt{\frac{1}{T}} \quad (57)$$

By substituting this value back, we have:

$$f(\bar{x}_T) - f(x^*) \leq B\rho \sqrt{\frac{1}{T}} \quad (58)$$



## 9 Theory Question

Based on the convexity property, we arrived at this particular form of expression by turning the coefficient into a sum of functions and incorporating them inside the function to make it a smaller set.

Indeed, let  $\eta$  be a vector, but the main intuition comes from the fact that we previously proved that for every random subgradients, we need to find  $\eta$  such that the function is decreasing in that direction. However, for this question, it is not necessary that all our steps be optimal. So, a reasonable approach is to not focus solely on the final solution and instead provide an effect from all the steps as an output. One such effect could be the arithmetic mean, which was examined in this question.

## 10 Theory Question

$$\sum_{t=1}^T \eta_t f(x_{(t)}) \geq \left( \sum_{t=1}^T \eta_t \right) f(x_{(t_T^*)}) \quad (59)$$

Using the above equation and also question 7, we have:

$$\sum_{t=1}^T \eta_t f(x_{(t)}) - a f(x^*) \leq \frac{1}{\gamma} B^{\gamma} + \frac{1}{\gamma} \rho^{\gamma} \sum_{t=1}^T \eta_t^{\gamma} \quad (60)$$

Deviding by  $a = \sum_{t=1}^T \eta_t$ , we have:

$$\rightarrow f(x_{(t_T^*)}) - f(x^*) \leq \frac{1}{\gamma} B^{\gamma} + \frac{\frac{1}{\gamma} \rho^{\gamma} \sum_{t=1}^T \eta_t^{\gamma}}{\sum_{t=1}^T \eta_t} \quad (61)$$

## 11 Theory Question

$$C = (x \in \mathbb{R}^n \mid Ax = b) \quad (62)$$

The goal is to find the closed-form expression for the projection of  $x$  onto  $C$ , denoted as  $\Pi_C(x)$ , which minimizes the squared Euclidean distance between  $y$  and  $x$ :

$$\Pi_C(x) = \arg \min_{y \in C} \|y - x\|^2 \quad (63)$$

To solve this optimization problem, we can utilize the concept of the Moore-Penrose pseudoinverse.

Formulating the Lagrangian function:

$$L(y, \lambda) = \|y - x\|^2 + \lambda^T (Ay - b) \quad (64)$$

where  $\lambda$  is a Lagrange multiplier vector associated with the equality constraint  $Ax = b$ . Taking the partial derivative of  $L$  with respect to  $y$  and setting it to zero:

$$\frac{\partial L}{\partial y} = 2(y - x) + A^T \lambda = 0 \quad (65)$$

This yields the following equation:

$$y = x - \frac{1}{2} A^T \lambda \quad (66)$$

Substituting this expression for  $y$  into the equality constraint  $Ax = b$ , we get:

$$Ax - \frac{1}{2} A A^T \lambda = b \quad (67)$$

To solve for  $\lambda$ , we can left-multiply the equation by  $A^\dagger$  (the left inverse of  $A$ ) to obtain:

$$A^\dagger Ax - \frac{1}{2} A^\dagger A A^T \lambda = A^\dagger b \quad (68)$$

Simplifying further:

$$\frac{1}{2} A^\dagger A A^T \lambda = x - A^\dagger b \quad (69)$$

Rearranging the equation:

$$\lambda = 2(x - A^\dagger b)^\dagger A^\dagger A (A^T)^\dagger \quad (70)$$

Multiplying both sides by 2 and solving for  $\lambda$ :

$$\lambda = \Upsilon (x - A^\dagger b)^\dagger A^\dagger A (A^T)^\dagger \quad (V1)$$

Finally, substituting the expression for  $\lambda$  back into  $y = x - 0.5A^T\lambda$ , we obtain the closed-form expression for  $\Pi_C(x)$ :

$$\Pi_C(x) = x - \frac{1}{2} A^T (A^\dagger A (A^T)^\dagger) (\Upsilon (x - A^\dagger b)) \quad (V2)$$

TA add the knowledge that  $AA^T$  is positive definite and  $m < n$  so the† will be gonna out and the final answer will be:

$$y = x - A^T (AA^T)^{-1} (Ax - b) \quad (V3)$$

This represents the projection of  $x$  onto the set  $C$  in closed form.

## 12 Theory Question

When the set  $C$  is defined as:

$$C = \{x \in \mathbb{R}^n \mid Ax \leq b\} \quad (V4)$$

where  $A$  is an  $m \times n$  matrix,  $x$  is an  $n$ -dimensional vector, and  $b$  is an  $m$ -dimensional vector, the closed-form expression for the projection of  $x$  onto  $C$ , denoted as  $\Pi_C(x)$ , can be found using the concept of linear programming and optimization.

The projection  $\Pi_C(x)$  can be defined as the solution to the following optimization problem:

$$\Pi_C(x) = \arg \min_{y \in C} \|y - x\|^\Upsilon \quad (V5)$$

To solve this problem, we can formulate it as a quadratic program with linear constraints. The objective function is to minimize the squared Euclidean distance between  $y$  and  $x$ :

$$\text{minimize: } \|y - x\|^\Upsilon \quad (V6)$$

subject to:  $Ax \leq b$

To find the closed-form expression, we can utilize the theory of quadratic programming and the KKT (Karush-Kuhn-Tucker) conditions.

The KKT conditions for this problem are as follows:

Stationarity condition:

$$\nabla(y - x) + A^\top v = 0 \quad (\text{V7})$$

Primal feasibility:

$$Ax - b \leq 0 \quad (\text{V8})$$

Dual feasibility:

$$v \geq 0 \quad (\text{V9})$$

Complementary slackness:

$$v^\top (Ax - b) = 0 \quad (\text{A10})$$

Here,  $v$  is a vector of Lagrange multipliers associated with the inequality constraints  $Ax \leq b$ .

By solving the KKT conditions, we can find the closed-form expression for  $\Pi_C(x)$ .

### 13 Theory Question

To find the center of the sphere and the intersection point of the line passing through the origin and a given point  $x$  with the sphere, we can proceed as follows:

Let the center of the sphere be denoted as  $c$ . If we consider the tangent line to the sphere at the intersection point  $x$ , we know that this line must pass through the center of the sphere. Using the right triangle formed by the line connecting  $x$ ,  $c$ , and the origin, we can prove that the line along  $x$  must coincide with the line connecting  $c$  and the origin.

Let's find the value of  $b$  relative to the center of the sphere. The intersection of this line segment with the sphere of radius  $b$  can be found by considering the line along  $(y : y = tx)$ .

$$y = \frac{x}{\|x\|_2} \cdot b \quad (A)$$

In this equation,  $c$  represents the center of the sphere,  $x$  is the given point, and  $b$  is the radius of the sphere. By multiplying the normalized vector  $x$  by the scalar  $b/\|x\|_2$ , we obtain the center of the sphere.

### 14 Theory Question

$$\begin{aligned} A : x^{t+1} &= \Pi_c(x^t - \eta_t v^t) = \operatorname{argmin}_{y \in C} \|y - (x^t - \eta_t v^t)\|_2 = \\ & \operatorname{argmin}_{y \in C} \|y - (x^t - \eta_t v^t)\|_2^2 = \operatorname{argmin}_{y \in C} \left\| y^T \cdot y - 2y^T(x^t - \eta_t v^t) + \|(x^t - \eta_t v^t)\|_2^2 \right\|_2 = \\ & \operatorname{argmin}_{y \in C} (\|y\|_2^2 - 2y^T(x^t - \eta_t v^t)) \\ B : x^{t+1} &= \operatorname{argmin}_{x \in C} (f(x^t) + \langle v^t, x - x^t \rangle + \frac{1}{2\eta_t} \|(x - x^t)\|) = \\ & \operatorname{argmin}_{x \in C} ((v^t)^T \cdot x + \frac{1}{2\eta_t} (\|x\|_2^2 - 2x^T x^t)) = \operatorname{argmin}_{x \in C} (\|x\|_2^2 - 2x^T(x^t - \eta_t v^t)) \\ & \xrightarrow{y=x} A \equiv B \end{aligned}$$

## 15 Theory Question

$$\begin{aligned}
 & \xrightarrow{\text{algorithm 2 theorems}} \frac{1}{2} \|x^{t+1} - x^*\|_2^2 \leq \frac{1}{2} \|x^t - x^*\|_2^2 - \eta_t (f(x^t) - f(x^*)) + \frac{\eta_t^2}{2} \|v^t\|_2^2 \\
 & \xrightarrow{\cdot \frac{1}{\eta_t}} \frac{1}{2\eta_t} \|x^{t+1} - x^*\|_2^2 \leq \frac{1}{2\eta_t} \|x^t - x^*\|_2^2 - (f(x^t) - f(x^*)) + \frac{\eta_t}{2} \|v^t\|_2^2 \\
 & \xrightarrow{\sum_{t=1}^T} \sum_{t=1}^T (f(x^t) - f(x^*)) \leq \frac{B^2}{2\eta_t} + \frac{\rho^2}{2} \sum_{t=1}^T \eta_t \\
 & \text{1.4,2.4 theorems}
 \end{aligned}$$

## 16 Theory Question

$$\begin{aligned}
 f(\bar{x}_t) &= f\left(\frac{x_1 + \dots + x_T}{T}\right) \xrightarrow[\text{jensen inequality}]{f \text{ is convex}} f(\bar{x}_t) \leq \frac{1}{T} \sum_{t=1}^T f(x^t) \\
 &\rightarrow f(\bar{x}_t) - f(x^*) \leq \left(\frac{1}{T} \sum_{t=1}^T f(x^t)\right) - f(x^*) \xrightarrow{15} \sum_{t=1}^T (f(x^t) - f(x^*)) \leq \frac{B^2}{2\eta_T} + \frac{\rho^2}{2} \sum_{t=1}^T \eta_t \\
 &\xrightarrow{* \frac{1}{T}} \left(\frac{1}{T} \sum_{t=1}^T f(x^t)\right) - f(x^*) \leq \frac{B^2}{2\eta_T T} + \frac{\rho^2}{2T} \sum_{t=1}^T \eta_t \rightarrow f(\bar{x}_t) - f(x^*) \leq \frac{B^2}{2\eta_T T} + \frac{\rho^2}{2T} \sum_{t=1}^T \eta_t \\
 &\xrightarrow{\eta_T = \frac{\alpha}{\sqrt{T}}} f(\bar{x}_t) - f(x^*) \leq \frac{B^2}{2\alpha\sqrt{T}} + \frac{\rho^2\alpha}{\sqrt{T}} \\
 &\xrightarrow{\text{opt. alpha}} \frac{\partial}{\partial \alpha} \left(\frac{B^2}{2\alpha\sqrt{T}} + \frac{\rho^2\alpha}{\sqrt{T}}\right) = 0 \rightarrow \left(-\frac{B^2}{2\alpha^2\sqrt{T}} + \frac{\rho^2}{\sqrt{T}}\right) = 0 \rightarrow \alpha = \frac{\sqrt{2}B}{\rho}
 \end{aligned}$$

## 17 Theory Question

با توجه به اینکه  $X$  از رابطه  $x^{t+1} = \Pi_c(x^t - \eta_t v^t)$  بدست می‌آید و از آنجاییکه در هر مرحله  $v^t$  از یک توزیع احتمالاتی بدست می‌آید و مقداری معین نیست،  $X$  بردار تصادفی است.

## 18 Theory Question

در بخش قبل به الگوریتم کاهش گرادیان مقید در حالتی که تابع هدف شامل متغیر تصادفی نیست رسیدیم. از آنجاییکه به دنبال روشی برای حل مسائل دارای متغیر تصادفی هستیم، منطقی است که به جای استفاده از زیرگرادیان‌های معلوم، به کمک اطلاعاتی که از متغیر تصادفی داریم (امید ریاضی) مقدار زیرگرادیان را حدس بزنیم.

## 19 Theory Question

$$E[v^t|x^t] \in \partial f(x^t) \rightarrow f(x^*) \geq f(x^t) + \langle x^* - x^t, E[v^t|x^t] \rangle$$

$$\xrightarrow{\eta_t > 0} -\eta_t((f(x^t) - f(x^*)) + \langle x^t - x^*, E[v^t|x^t] \rangle) > 0$$

حال کافی است ثابت کنیم:

$$\left(\frac{1}{2}\right) \|x^{t+1} - x^*\|_2^2 \leq \frac{1}{2} \|x^t - x^*\|_2^2 + \frac{\eta_t^2}{2} \|v^t\|_2^2 - \eta_t \langle v^t, x^t - x^* \rangle = \frac{1}{2} \|(x^t - \eta_t v^t) - x^*\|_2^2$$

برهان خلف: اگر رابطه بالا برقرار نباشد، با توجه به رابطه  $x^{t+1} = \Pi_c(x^t - \eta_t v^t)$  می‌توان گفت  $x^*$  نزدیک‌ترین نقطه به  $x^t - \eta_t v^t$  است پس  $x^{t+1} = x^*$  و در نتیجه باید داشته باشیم:

$$\|(x^t - \eta_t v^t) - x^*\|_2^2 < 0$$

از آنجاییکه نامساوی فوق ناممکن است فرض خلف باطل و حکم برقرار است.

## 20 Theory Question

$$\begin{aligned} E[\langle \varepsilon_t, X^t - x^* \rangle] &= E[(V^t)^T X^t - (V^t)^T x^* - E[V^t|X^t]^T X^t + E[V^t|X^t]^T x^*] = \\ &= E[V^t X^t] - E[(V^t)^T] x^* - E[E[V^t|X^t]^T X^t] + E[V^t]^T x^* = \\ &= E[(V - E[V])^T (X - E[X])] - E[(E[V|X] - E[E[V|X]])^T (X - E[X])] = \\ &= E[(V - E[V])^T (X - E[X])] - E[(E[V|X] - E[V])^T (X - E[X])] = 0 \end{aligned}$$

## 21 Theory Question

$$\begin{aligned} \frac{1}{2} \|X^{t+1} - x^*\|_2^2 &\leq \frac{1}{2} \|X^t - x^*\|_2^2 - \eta_t (f(X^t) - f(x^*)) + \frac{\eta_t^2}{2} \|V^t\|_2^2 - \eta_t \langle \varepsilon_t, X - x^* \rangle \\ \xrightarrow{\cdot \frac{1}{\eta_t}} \frac{1}{2\eta_t} \|X^{t+1} - x^*\|_2^2 &\leq \frac{1}{2\eta_t} \|X^t - x^*\|_2^2 - (f(X^t) - f(x^*)) + \frac{\eta_t}{2} \|V^t\|_2^2 - \langle \varepsilon_t, X - x^* \rangle \\ \xrightarrow[\substack{1.4, 2.4 \text{ theorems}}]{\substack{\sum_{t=1}^T, Q_{16} \text{ proof}}} f(\bar{X}^t) - f(x^*) &\leq \frac{B^2}{2T\eta_T} + \frac{\rho^2}{2T} \sum_{t=1}^T \eta_t - \frac{1}{T} \sum_{t=1}^T \eta_t \langle \varepsilon_t, X_t - x^* \rangle \\ \xrightarrow[\substack{Q_{20}}]{E[\cdot]} E[f(\bar{X}^t)] - f(x^*) &\leq \frac{B^2}{2T\eta_T} + \frac{\rho^2}{2T} \sum_{t=1}^T \eta_t \end{aligned}$$

## 22 Theory Question

## 23 Theory Question

## 24 Theory Question

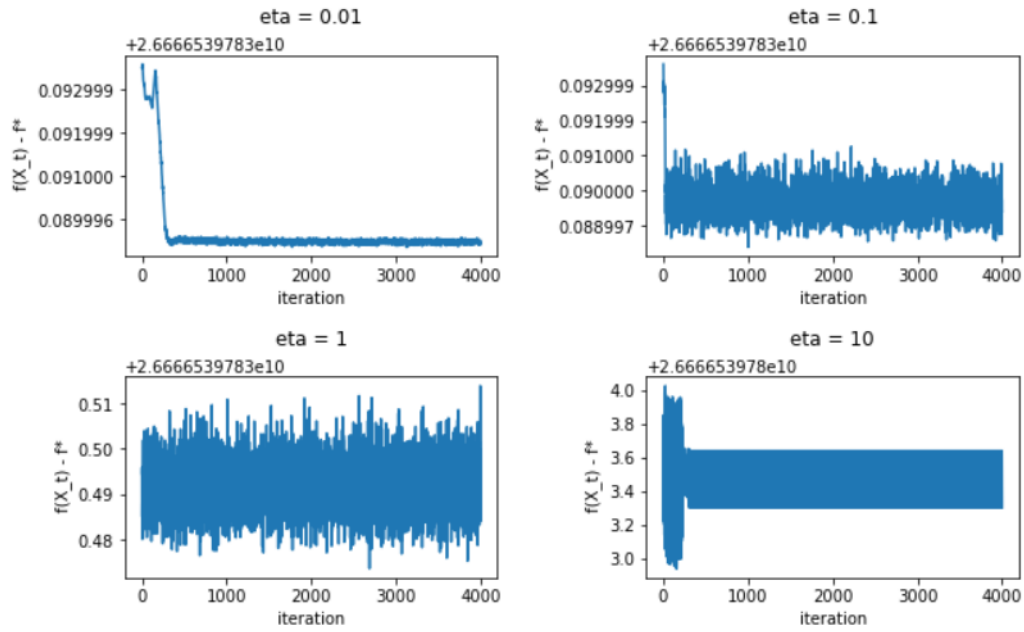
$$v \in \partial F(x, z) \xrightarrow[\text{course slides}]{p(x) > 0} p(x)v \in \partial(p(x)F(x))$$

$\xrightarrow{f} E[V] \in \partial E[F] = \partial f(X) \rightarrow v \text{ is random subgradient.}$



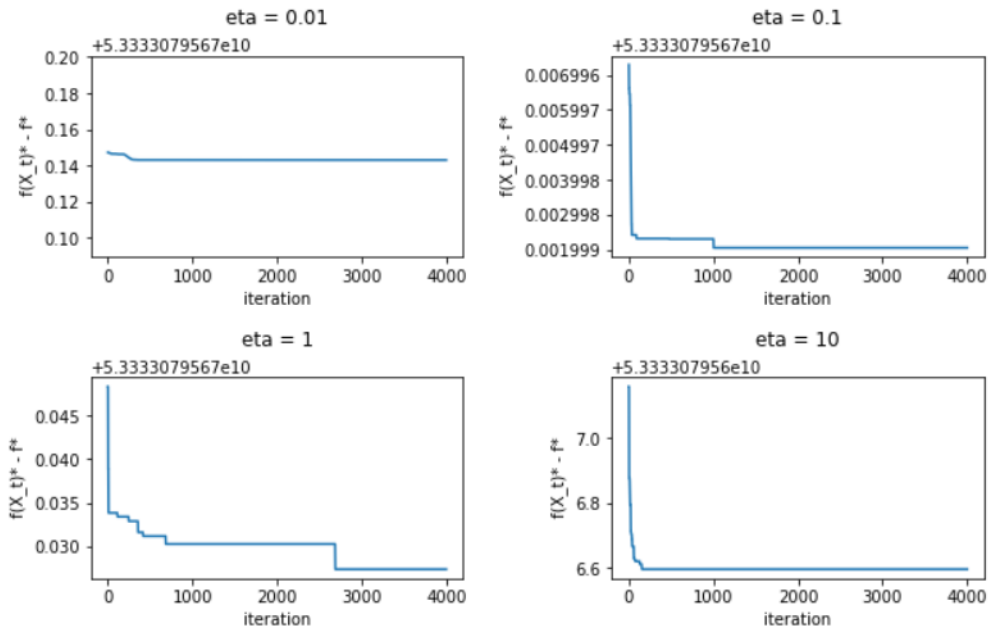
# Simulation

1.



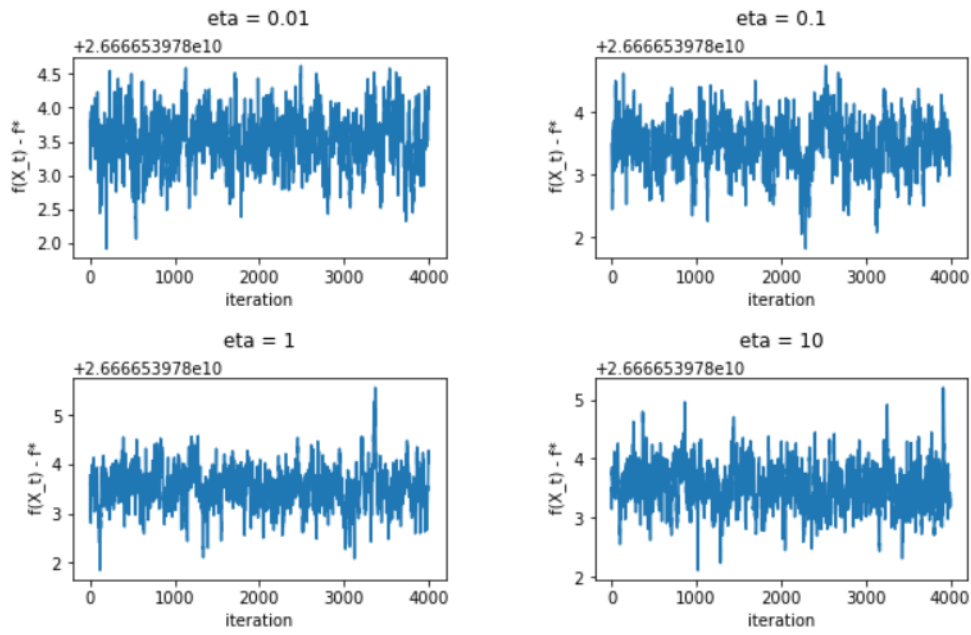
for small  $\eta$  error converge to zero.

2.



error converge to zero faster than part 1.

3.



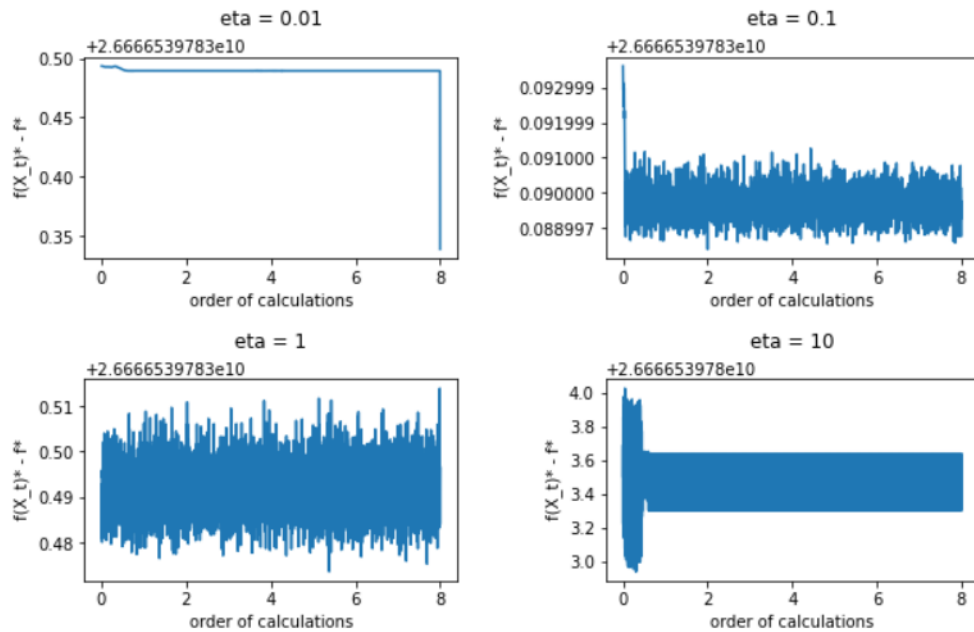
it doesn't converge while in part 1 the error converges to zero.

4.

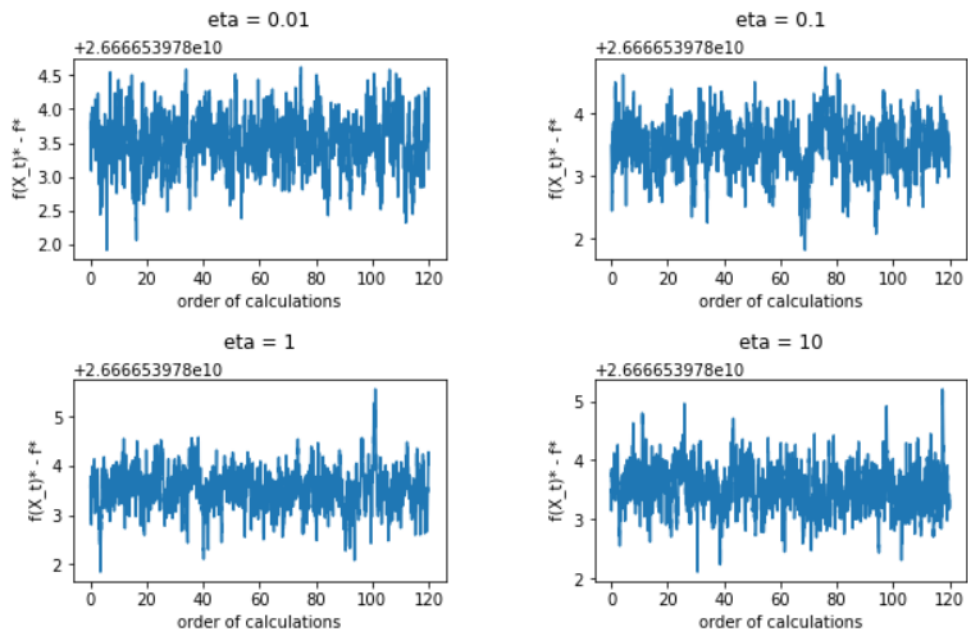
Elapsed time: 0.0020360946655273438 seconds

Elapsed time: 0.029249191284179688 seconds

numbers above are the times spent for an iteration of each algorithm.



error over calculation order plot for algorithm 2.



error over calculation order plot for algorithm 4.  
 based on plots above, it is obvious that algorithm 4 takes much more calculation power compared to algorithm 2.  
 rest of the questions are explained in notebook.