

In the name of God

Convex optimization project

Ardalan Gerami 99102112

Mahyar Ghazanfari 98102057



Q1.

we know :

$$X^{(t+1)} = X^{(t)} - \eta \nabla f(X^{(t)}), X^{(1)} = 0$$

and also :

$$\nabla f(X^{(t)}) = AX^{(t)} - b$$

so we can write \rightarrow

$$X^{(t+1)} = X^{(t)} - \eta AX^{(t)} + \eta b = (I - \eta A)X^{(t)} + \eta b$$

now from :

$$X^* = \underset{X}{\operatorname{argmin}} \{f(X)\} \rightarrow AX^* - b = 0$$

$$\rightarrow X^{(t+1)} - X^* = X^{(t)} - \eta(AX^{(t)} - b) - X^* = X^{(t)} - \eta(AX^{(t)} - b) - X^* + \eta(AX^* - b)$$

$$\rightarrow X^{(t+1)} - X^* = (I - \eta A)(X^{(t)} - X^*) = (I - \eta A) \dots (I - \eta A)(X^{(1)} - X^*)$$

$$\rightarrow X^{(t+1)} - X^* = -(I - \eta A)^t X^*, \|X^{(t+1)} - X^*\| \leq \|(I - \eta A)^t\| \|X^*\|$$

$$\rightarrow \|X^{(t+1)} - X^*\| \leq \|(I - \eta A)\|^t \|X^*\|$$

on the other hand if we want to have $\lim_{T \rightarrow \infty} \|X^T - X^*\| = 0$

we should *consider f as a M_{smooth} and $m_{strongly}$ convex :*

$$mI \leq \nabla^2 f(X) \leq MI \rightarrow \frac{m}{M}I \leq \frac{A}{M} \leq I$$

also :

$$\left\| I - \frac{A}{\lambda_{\max}(A)} \right\|_2 = 1 - \frac{\lambda_{\min}(A)}{\lambda_{\max}(A)}$$

$$\rightarrow \|X^{(t+1)} - X^*\| \leq \left(1 - \frac{m}{M}\right)^t \|X^*\| \leq e^{-\frac{m}{M}t} \|X^*\|, \text{ now if } t \rightarrow \infty : \|X^{(t+1)} - X^*\| \leq 0$$

so

$$\lim_{T \rightarrow \infty} \|X^T - X^*\| = 0$$

Q2.

we want to prove

$$X^{(t+1)} = X^{(t)} - \eta \nabla f(X^{(t)}) = \arg \min_X \{ f(X^{(t)}) + \langle \nabla f(X^{(t)}), X - X^{(t)} \rangle + \frac{1}{2\eta} \|X - X^{(t)}\|_2^2 \}$$

let's calculate the gradient of

$$h(f) = f(X^{(t)}) + \langle \nabla f(X^{(t)}), X - X^{(t)} \rangle + \frac{1}{2\eta} \|X - X^{(t)}\|_2^2$$

$$\rightarrow \nabla_X h(f) = \nabla_X \left(f(X^{(t)}) + \langle \nabla f(X^{(t)}), X - X^{(t)} \rangle + \frac{1}{2\eta} \|X - X^{(t)}\|_2^2 \right)$$

$$\nabla_X h(f) = \nabla_X \langle \nabla f(X^{(t)}), X - X^{(t)} \rangle + \nabla_X \frac{1}{2\eta} \|X - X^{(t)}\|_2^2$$

$$\nabla_X h(f) = \nabla f(X^{(t)}) + \frac{1}{2\eta} (2X - 2X^{(t)}) = 0$$

$$\nabla f(X^{(t)}) + \frac{(X - X^{(t)})}{\eta} = 0 \rightarrow X = X^{(t)} - \eta \nabla f(X^{(t)}) \text{ and } X^{(t+1)} = X$$

so

$$X^{(t+1)} = \arg \min_X \{ f(X^{(t)}) + \langle \nabla f(X^{(t)}), X - X^{(t)} \rangle + \frac{1}{2\eta} \|X - X^{(t)}\|_2^2 \}$$

Q3.

$$\|X^{(t+1)} - X^*\|_2^2 = \|X^{(t)} - \eta_t v^{(t)} - X^*\|_2^2 \leq \|X^{(t)} - X^*\|_2^2$$

above all we know:

$$\|X^{(t)} - \eta_t v^{(t)} - X^*\|_2^2 = \|X^{(t)} - X^*\|_2^2 + \eta_t^2 \|v^{(t)}\|_2^2 - 2\langle X^{(t)} - X^*, \eta_t v^{(t)} \rangle$$

on the other hand :

$$\|X^{(t)} - X^*\|_2^2 + \eta_t^2 \|v^{(t)}\|_2^2 - 2\langle X^{(t)} - X^*, \eta_t v^{(t)} \rangle \leq \|X^{(t)} - X^*\|_2^2 - \eta_t^2 \|v^{(t)}\|_2^2$$

and:

$$\|X^{(t)} - X^*\|_2^2 - \eta_t^2 \|v^{(t)}\|_2^2 \leq \|X^{(t)} - X^*\|_2^2$$

so :

$$\|X^{(t+1)} - X^*\|_2^2 \leq \|X^{(t)} - X^*\|_2^2$$

Q4.

due to the definition of sub - gradient such v is sub - gradient of f in point $x, y \in \mathbb{R}^n$ so:

$$f(y) \geq f(x) + \langle v, y - x \rangle$$

and f is ρ - lipshitz :

$$|f(y) - f(x)| \leq \rho \|y - x\|_2$$

go on :

$$|f(y) - f(x)| \geq \|v\|_2 \|y - x\|_2$$

so :

$$\|v\|_2 \leq \rho$$

Q5.

we know(I) :

$$\|X^{(t+1)} - X^*\|_2^2 = \|X^{(t)} - \eta_t v^{(t)} - X^*\|_2^2 = \|X^{(t)} - X^*\|_2^2 + \eta_t^2 \|v^{(t)}\|_2^2 - 2\eta_t \langle X^{(t)} - X^*, v^{(t)} \rangle$$

on the other hand $X^* = \operatorname{argmin} f(x)$ so we would have :

$$f(X^{(t)}) \geq f(X^*)$$

and :

$$|f(X^{(t)}) - f(X^*)| \leq \rho \|X^{(t)} - X^*\|_2, \|v^{(t)}\|_2 \leq \rho$$

so :

$$f(X^{(t)}) - f(X^*) \leq \rho \|X^{(t)} - X^*\|_2 \rightarrow \text{multiply } 2\eta_t \|v^{(t)}\|_2, \text{ we have:}$$

$$2\eta_t \|v^{(t)}\|_2 (f(X^{(t)}) - f(X^*)) \leq 2\eta_t \|v^{(t)}\|_2 \rho \|X^{(t)} - X^*\|_2$$

on the other hand we said $\|v^{(t)}\|_2 \leq \rho$, so :

$$2\eta_t (f(X^{(t)}) - f(X^*)) \leq 2\eta_t \|v^{(t)}\|_2 \|X^{(t)} - X^*\|_2$$

now we try to build (I):

$$-2\eta_t (f(X^{(t)}) - f(X^*)) \geq -2\eta_t \|v^{(t)}\|_2 \|X^{(t)} - X^*\|_2$$

\rightarrow

$$\|X^{(t)} - X^*\|_2^2 + \eta_t^2 \|v^{(t)}\|_2^2 - 2\eta_t (f(X^{(t)}) - f(X^*))$$

$$\geq \|X^{(t)} - X^*\|_2^2 + \eta_t^2 \|v^{(t)}\|_2^2 - 2\eta_t \|v^{(t)}\|_2 \|X^{(t)} - X^*\|_2$$

so :

$$\|X^{(t+1)} - X^*\|_2^2 \leq \|X^{(t)} - X^*\|_2^2 + \eta_t^2 \|v^{(t)}\|_2^2 - 2\eta_t (f(X^{(t)}) - f(X^*))$$

$$\rightarrow \times \frac{1}{2} :$$

$$\frac{1}{2} \|X^{(t+1)} - X^*\|_2^2 \leq \frac{1}{2} \|X^{(t)} - X^*\|_2^2 - \eta_t (f(X^{(t)}) - f(X^*)) + \frac{\eta_t^2}{2} \|v^{(t)}\|_2^2$$

Q6.

first consider $X_0 = 0 = X^{(1)}$, also We consider the assumptions of the question :

$$\sum_{t=1}^T \eta_t (f(X^{(t)}) - f(X^*)) \leq \frac{1}{2} B^2 + \frac{1}{2} \rho^2 \sum_{t=1}^T \eta_t^2$$

on the other hand $\|v^{(t)}\|_2 \leq \rho$ and $\|X^*\|_2^2 \leq B^2$ and from part 3:

$$\frac{1}{2} \|X^{(t+1)} - X^*\|_2^2 \leq \frac{1}{2} \|X^{(t)} - X^*\|_2^2 - \eta_t (f(X^{(t)}) - f(X^*)) + \frac{\eta_t^2}{2} \|v^{(t)}\|_2^2$$

$$\rightarrow \sum_{t=1}^T (\text{each side of inequality}):$$

$$\frac{1}{2} \sum_{t=1}^T \|X^{(t+1)} - X^*\|_2^2 \leq \frac{1}{2} \sum_{t=1}^T \|X^{(t)} - X^*\|_2^2 - \sum_{t=1}^T \eta_t (f(X^{(t)}) - f(X^*)) + \frac{1}{2} \sum_{t=1}^T \eta_t^2 \|v^{(t)}\|_2^2$$

subtraction :

$$\frac{1}{2} (\|X^{(T)} - X^*\|_2^2) \leq - \sum_{t=1}^T \eta_t (f(X^{(t)}) - f(X^*)) + \frac{1}{2} \sum_{t=1}^T \eta_t^2 \|v^{(t)}\|_2^2$$

$$\rightarrow \|X^{(T)} - X^*\|_2^2 \leq \|X^{(1)} - X^*\|_2^2 = \|X^*\|_2^2 \leq B^2, \|v^{(t)}\|_2 \leq \rho, \text{ so:}$$

$$\sum_{t=1}^T \eta_t (f(X^{(t)}) - f(X^*)) \leq \frac{1}{2} \rho^2 \sum_{t=1}^T \eta_t^2 + \frac{1}{2} B^2$$

Q7.

let's divide $\sum_{t=1}^T \eta_t (f(X^{(t)}) - f(X^*))$ by $\sum_{t=1}^T \eta_t^2$:

$$\frac{\sum_{t=1}^T \eta_t (f(X^{(t)}) - f(X^*))}{\sum_{t=1}^T \eta_t^2}$$

if we consider the form of the $X_T^{bar} = \frac{\sum_{t=1}^T \eta_t X^{(t)}}{\sum_{t=1}^T \eta_t}$ we conclude :

$$\frac{\sum_{t=1}^T \eta_t (f(X^{(t)}) - f(X^*))}{\sum_{t=1}^T \eta_t^2} = f(X_T^{bar}) - f(X^*)$$

due to Q6 and divide each side of inequality by $\sum_{t=1}^T \eta_t^2$ we have :

$$\frac{\sum_{t=1}^T \eta_t (f(X^{(t)}) - f(X^*))}{\sum_{t=1}^T \eta_t^2} = f(X_T^{bar}) - f(X^*) \leq \frac{B^2 + \rho^2 \sum_{t=1}^T \eta_t^2}{2 \sum_{t=1}^T \eta_t}$$

so:

$$f(X_T^{bar}) - f(X^*) \leq \frac{B^2 + \rho^2 \sum_{t=1}^T \eta_t^2}{2 \sum_{t=1}^T \eta_t}$$

Q8.

*** we should consider $t = 0$ to T ***

let $\eta = \eta_t$:

$$f(X_T^{bar}) - f(X^*) \leq \frac{B^2 + (T+1)\rho^2\eta^2}{2(T+1)\eta} = \frac{1}{2} \left(\frac{B^2}{(T+1)\eta} + \eta\rho^2 \right)$$

for best η :

$$\operatorname{argmin} \frac{1}{2} \left(\frac{B^2}{(T+1)\eta} + \eta\rho^2 \right) \rightarrow \frac{\rho^2}{2} - \frac{B^2}{2(T+1)\eta^2} = 0 \rightarrow \eta_* = \frac{B}{\rho\sqrt{T+1}}$$

best bound:

$$\frac{1}{2} \left(\frac{B^2}{(T+1)\eta} + \eta\rho^2 \right), \eta_* = \frac{B}{\rho\sqrt{T+1}} \rightarrow \frac{\rho B}{\sqrt{T+1}}$$

Q9.

با توجه به اینکه در *Sub-Gradient Descent* لزوماً در هر مرحله کاهش انجام نمی‌گیرد و بعضاً نقاطی وجود دارند که افزایش می‌یابند می‌توانیم برای این نقاط از میانگین تمامی نقاط در الگوریتم استفاده کنیم و با توجه به اینکه در سوال بالا کران نهایی بر حسب $f(X_T^{bar})$ به دست آمد برای مقدار تابع در میانگین کران بالا وجود دارد ولی برای آن نقطه به روزرسانی شده که افزایشی است پیدا کردن کران بالا مقدور نیست.

Q10.

$$f(X_T^{bar}) - f(X^*) \leq \frac{B^2 + \rho^2 \sum_{t=1}^T \eta_t^2}{2 \sum_{t=1}^T \eta_t} \text{ (from question 7):}$$

$$\frac{\sum_{t=1}^T \eta_t (f(X^{(t)}))}{\sum_{t=1}^T \eta_t^2} - f(X^*) \leq \frac{B^2 + \rho^2 \sum_{t=1}^T \eta_t^2}{2 \sum_{t=1}^T \eta_t}$$

on the other hand $f(X^{(t^*)}) \leq f(X^{(t)})$ so:

$$\frac{\sum_{t=1}^T \eta_t f(X^{(t^*)})}{\sum_{t=1}^T \eta_t^2} \geq f(X^{(t^*)})$$

so:

$$f(X^{(t^*)}) - f(X^*) \leq \frac{B^2 + \rho^2 \sum_{t=0 \text{ or } 1}^T \eta_t^2}{2 \sum_{t=1 \text{ or } 0}^T \eta_t}$$

Q11.

we can write the $\Pi_c = \underset{y}{\operatorname{argmin}} \{ \|y - x\|_2 \}$ as a convex optimization problem:

$$\begin{aligned} & \text{minimize } \|y - x\| \\ & \text{subject to } Ay = b \end{aligned}$$

solve :

by lagrangian coefs we have :

$$\mathcal{L}(x, y, v) = x^T x + y^T y - 2x^T y + v^T (Ay - b)$$

calculate $\frac{\partial \mathcal{L}}{\partial y}$:

$$\frac{\partial \mathcal{L}}{\partial y} = 2y - 2x + A^T v = 0 \rightarrow y_* = \frac{1}{2}(2x - A^T v) \rightarrow \times A, Ay_* = b, \text{ so:}$$

$$Ay_* = Ax - \frac{1}{2}AA^T v \rightarrow \frac{1}{2}AA^T v = Ax - b$$

$$v = 2(AA^T)^{-1}(Ax - b)$$

now we calculate y :

$$y_* = x - 2A^T(AA^T)^{-1}(Ax - b)$$

so:

$$\Pi_c = \left\| 2A^T(AA^T)^{-1}(Ax - b) \right\|_2$$

Q12.

we can write the $\Pi_c = \underset{y}{\operatorname{argmin}} \{ \|y - x\|_2 \}$ as a convex optimization problem:

$$\text{minimize } \|y - x^2\|$$

$$\text{subject to } Ay \leq b$$

solve :

by lagrangian coefs we have :

$$\mathcal{L}(y, \lambda) = \|y - x^2\| + \lambda^T(Ay - b)$$

$$\text{calculate } \frac{\partial \mathcal{L}}{\partial y} :$$

$$\frac{\partial \mathcal{L}}{\partial y} = 2y - 2x + A^T \lambda = 0$$

$$KKT: \begin{cases} 2y - 2x + A^T \lambda = 0 \\ Ay \leq b \\ \lambda_i(A_i^T \lambda - b_i) = 0 \end{cases}$$

Q13.

C determine the set of inner points of a sphere and we can say:

The closest point of the set to the desired point is in the line between this point and the center of the sphere (on the sphere). Therefore, y is in the direction of x and at a distance b from the center of the sphere.

خلاصه: y در راستای x در فاصله b از مرکز کره است.

Q14.

due to algorithm 3 we have :

$$X^{(t+1)} = \Pi_c(X^{(t)} - \eta_t v^{(t)})$$

we want to prove that above $X^{(t+1)}$ is same as below eq.:

$$X^{(t+1)} = \underset{X}{\operatorname{argmin}} \{f(X^{(t)}) + \langle v^{(t)}, X - X^{(t)} \rangle + \frac{1}{2\eta_t} \|X - X^{(t)}\|_2^2\}$$

go on:

$$\begin{aligned} X^{(t+1)} &= \Pi_c(X^{(t)} - \eta_t v^{(t)}) = \underset{y}{\operatorname{argmin}} \|y - X^{(t)} + v^{(t)}\|_2 \\ &= \underset{y}{\operatorname{argmin}} \left(y^T y - 2y^T (X^{(t)} - \eta_t v^{(t)}) + \|X^{(t)} - \eta_t v^{(t)}\|_2^2 \right) \\ &= \underset{y}{\operatorname{argmin}} (\|y\|_2^2 - 2y^T (X^{(t)} - \eta_t v^{(t)})) \end{aligned}$$

from eq.:

$$X^{(t+1)} = \underset{X}{\operatorname{argmin}} \{f(X^{(t)}) + \langle v^{(t)}, X - X^{(t)} \rangle + \frac{1}{2\eta_t} \|X - X^{(t)}\|_2^2\}$$

$$X^{(t+1)} = \underset{X}{\operatorname{argmin}} \{v^{(t)} \cdot x + \frac{1}{2\eta_t} (\|x\|_2^2 - 2x^T x^{(t)})\}$$

$$X^{(t+1)} = \underset{y}{\operatorname{argmin}} (\|x\|_2^2 - 2x^T (X^{(t)} - \eta_t v^{(t)}))$$

as we see Constrained Sub – Gradient Descent method will leads to same answer.

Q15.

as we know from th. 6 we can write :

$$\frac{1}{2} \|X^{(t+1)} - X^*\|_2^2 \leq \frac{1}{2} \|X^{(t)} - X^*\|_2^2 - \eta_t (f(X^{(t)}) - f(X^*)) + \frac{\eta_t^2}{2} \|v^{(t)}\|_2^2$$

by multiplying $\frac{1}{\eta_t}$ and $\sum_{t=1}^T$ (each side of inequality)

$$\rightarrow \frac{1}{2\eta_t} (\|X^{(T)} - X^*\|_2^2) \leq - \sum_{t=1}^T \eta_t (f(X^{(t)}) - f(X^*)) + \frac{1}{2} \sum_{t=1}^T \eta_t^2 \|v^{(t)}\|_2^2$$

\rightarrow

on the other hand $\|v^{(t)}\|_2 \leq \rho$ and $\|X - X^*\|_2^2 \leq B^2$:

$$\sum_{t=1}^T (f(X^{(t)}) - f(X^*)) \leq \frac{1}{2} \rho^2 \sum_{t=1}^T \eta_t + \frac{1}{2\eta_T} B^2$$

Q16.

$$\text{let } \eta_t = \frac{\alpha}{\sqrt{t}} \text{ and } X_T^{\text{bar}} = \frac{1}{T} \sum_{t=1}^T X^{(t)} :$$

start algorithm 3 with $X_0 \in C$:

from jensen inequality :

$$f(X_T^{\text{bar}}) \leq \frac{1}{T} \sum_{t=1}^T f(X^{(t)}) , f(X_T^{\text{bar}}) - f(X^*) \leq \frac{1}{T} \sum_{t=1}^T f(X^{(t)}) - f(X^*)$$

from last part we know :

$$\frac{1}{T} \sum_{t=1}^T f(X^{(t)}) - f(X^*) \leq \frac{1}{2} \rho^2 \sum_{t=1}^T \eta_t + \frac{1}{2\eta_T} B^2 \rightarrow$$

$$\sum_{t=1}^T f(X^{(t)}) - f(X^*) \leq \frac{1}{2T} \rho^2 \sum_{t=1}^T \eta_t + \frac{1}{2T\eta_T} B^2$$

$$f(X_T^{\text{bar}}) - f(X^*) \leq \frac{\rho^2 \alpha}{\sqrt{T}} + \frac{1}{2\alpha\sqrt{T}} B^2$$

now optimal α :

$$\alpha_* = \operatorname{argmin} \frac{\rho^2 \alpha}{\sqrt{T}} + \frac{1}{2\alpha\sqrt{T}} B^2 \rightarrow \alpha_* = \frac{B\sqrt{2}}{\rho}$$

Q17.

می‌دانیم که $v^{(t)}$ ناشی از یک توزیع احتمالاتی است و همچنین $X^{(t)}$ در هر لحظه 0 تا $T+1$ از رابطه

$$X^{(t+1)} = \Pi_c(X^{(t)} - \eta_t v^{(t)})$$

بدست می‌آید بنابراین می‌توان استدلال نمود X برداری تصادفی است و منشأ این خاصیت تصادفی همین بیانات بالا است.

Q18.

تا به الان هیچگاه تابع هدف متغیر تصادفی را شامل نمیشد. در این مسئله ما *gradient-discent* را برای متغیرهای تصادفی که توزیع تصادفی دارند تعریف کردیم و باید مجموعه *subgradient* را باید با استفاده از امید ریاضی داده‌ها و بردارها بدست بیاوریم.

Q19.

we know $\mathcal{E}_t = V^{(t)} - E[V^{(t)}|X^{(t)}]$ and we want to prove :

$$\frac{1}{2} \|X^{(t+1)} - X^*\|_2^2 \leq \frac{1}{2} \|X^{(t)} - X^*\|_2^2 - \eta_t (f(X^{(t)}) - f(X^*)) + \frac{\eta_t^2}{2} \|v^{(t)}\|_2^2 - \eta_t \langle \mathcal{E}_t, X^{(t)} - X^* \rangle$$

go on :

due to question we know $E[V^{(t)}|X^{(t)}] \in \partial f(X^{(t)})$ so :

$$f(X^*) \geq f(X^{(t)}) + \langle X^{(t)} - X^*, E[V^{(t)}|X^{(t)}] \rangle \rightarrow \text{multiply } \eta_t:$$

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

we have to prove equivalent :

$$\frac{1}{2} \|X^{(t+1)} - X^*\|_2^2 \leq \frac{1}{2} \|X^{(t)} - \eta_t v^{(t)} - X^*\|_2^2 \quad (I)$$

as we know $X^{(t+1)} = \Pi_c(X^{(t)} - \eta_t v^{(t)})$ and if (I) doesn't be valid then $X^{(t+1)} = X^*$

$$\text{and } \frac{1}{2} \|X^{(t+1)} - X^*\|_2^2 \leq \frac{1}{2} \|X^{(t)} - \eta_t v^{(t)} - X^*\|_2^2 \leq 0 \text{ which is in not valid so:}$$

$$\frac{1}{2} \|X^{(t+1)} - X^*\|_2^2 \leq \frac{1}{2} \|X^{(t)} - X^*\|_2^2 - \eta_t (f(X^{(t)}) - f(X^*)) + \frac{\eta_t^2}{2} \|v^{(t)}\|_2^2 - \eta_t \langle \mathcal{E}_t, X^{(t)} - X^* \rangle$$

Q20.

$$\text{show } E[\langle \mathcal{E}_t, X^{(t)} - X^* \rangle] = 0$$

$$E[\langle \mathcal{E}_t, X^{(t)} - X^* \rangle] = E[\langle V^{(t)} - E[V^{(t)}|X^{(t)}], X^{(t)} - X^* \rangle] = E[\langle V^{(t)}, X^{(t)} - X^* \rangle] - E[\langle E[V^{(t)}|X^{(t)}], X^{(t)} - X^* \rangle]$$

$$E[\langle V^{(t)}, X^{(t)} - X^* \rangle] - E[\langle E[V^{(t)}|X^{(t)}], X^{(t)} - X^* \rangle] = E[(X^{(t)} - X^*)^T V^{(t)}] - E[(X^{(t)} - X^*)^T E[V^{(t)}|X^{(t)}]]$$

$$E[(X^{(t)} - X^*)^T V^{(t)}] - E[(X^{(t)} - X^*)^T E[V^{(t)}|X^{(t)}]] = E[(X^{(t)} - X^*)^T V^{(t)}] - E[E[(X^{(t)} - X^*)^T V^{(t)}|X^{(t)}]]$$

so :

$$E[(X^{(t)} - X^*)^T V^{(t)}] - E[(X^{(t)} - X^*)^T V^{(t)}] = 0 \rightarrow$$

$$E[\langle \mathcal{E}_t, X^{(t)} - X^* \rangle] = 0$$

Q21.

$$E[f(X_T^{bar})] - f(X^*) = E\left[f\left(\frac{1}{T} \sum_{t=1}^T X^{(t)}\right)\right] - f(X^*) \leq \frac{1}{T} \sum_{t=1}^T E[f(X^{(t)})] - f(X^*) = \frac{1}{T} \sum_{t=1}^T E[f(X^{(t)}) - f(X^*)]$$

$$\text{so : } E[f(X_T^{bar})] - f(X^*) \leq \frac{1}{T} \sum_{t=1}^T E[f(X^{(t)}) - f(X^*)]$$

from question 19 , 20 :

$$f(X^{(t)}) - f(X^*) \leq \frac{1}{2\eta_t} \|X^{(t)} - X^*\|_2^2 - \frac{1}{2\eta_t} \|X^{(t+1)} - X^*\|_2^2 + \frac{\eta_t}{2} \|v^{(t)}\|_2^2 - \langle \mathcal{E}_t, X^{(t)} - X^* \rangle$$

$$E[f(X^{(t)}) - f(X^*)] \leq \frac{1}{2\eta_t} E[\|X^{(t)} - X^*\|_2^2] - \frac{1}{2\eta_t} E[\|X^{(t+1)} - X^*\|_2^2] + \frac{\eta_t}{2} E[\|v^{(t)}\|_2^2]$$

$$\rightarrow \sum_{t=1}^T (\text{each side of inequality}):$$

$$\rightarrow \sum_{t=1}^T E[f(X^{(t)}) - f(X^*)] \leq \frac{1}{2\eta_t} E[\|X^{(1)} - X^*\|_2^2] - \frac{1}{2\eta_{T+1}} E[\|X^{(T+1)} - X^*\|_2^2] + \sum_{t=1}^T \frac{\eta_t}{2} E[\|v^{(t)}\|_2^2]$$

$$\text{from } E[\|v\|_2^2] \leq \rho^2, \|X - X^*\|_2^2 \leq B^2:$$

$$\sum_{t=1}^T E[f(X^{(t)}) - f(X^*)] \leq \frac{1}{2\eta_t} B^2 - \frac{1}{2\eta_{T+1}} E[\|X^{(T+1)} - X^*\|_2^2] + \sum_{t=1}^T \frac{\eta_t}{2} \rho^2$$

go on:

$$\sum_{t=1}^T E[f(X^{(t)}) - f(X^*)] \leq \frac{1}{2\eta_t} B^2 + \sum_{t=1}^T \frac{\eta_t}{2} \rho^2$$

$$\mathbb{E}[f(X_T^{bar})] - f(X^*) \leq \frac{1}{T} \sum_{t=1}^T \mathbb{E}[f(X^{(t)}) - f(X^*)] \leq \frac{1}{2\eta_t} B^2 + \sum_{t=1}^T \frac{\eta_t}{2} \rho^2$$

so :

$$\mathbb{E}[f(X_T^{bar})] - f(X^*) \leq \frac{1}{2\eta_t} B^2 + \sum_{t=1}^T \frac{\eta_t}{2} \rho^2$$

Q22.

let $\eta_t = \frac{B}{\rho\sqrt{t}}$ and from Q21 :

$$\mathbb{E}[f(X_T^{bar})] - f(X^*) \leq \frac{1}{2T\eta_t} B^2 + \frac{1}{2T} \rho^2 \sum_{t=1}^T \eta_t = \frac{B\rho}{2\sqrt{T}} + \frac{B\rho}{2T} \sum_{t=1}^T \frac{1}{\sqrt{t}}$$

$$\frac{B\rho}{2\sqrt{T}} + \frac{B\rho}{2T} \sum_{t=1}^T \frac{1}{\sqrt{t}} \leq \frac{3B\rho}{2\sqrt{T}} \rightarrow \sum_{t=1}^T \frac{1}{\sqrt{t}} \leq 2\sqrt{T}$$

by induction and let $T = 1 : 1 < 2$ which is correct

for $T + 1$ we should show $2\sqrt{T} \leq 2\sqrt{T+1} - \frac{1}{\sqrt{T+1}}$ if we square each side :

$$\frac{1}{T+1} \geq 0 \rightarrow T \geq 1 \text{ so by induction proved}$$

Q23.

Q24.

$\mathbb{E}[\cdot]$ Mathematical hope is a linear operator so:

$$\forall \partial_X F(X, Z) \rightarrow$$

$$\mathbb{E}[\partial_X F(X, Z)] = \partial \mathbb{E}_Z[F(X, Z)] = \partial f(X)$$

so V is a Random subgradient

Simulation questions:

Simulation 1:

First of all, It is better to rewrite function F as below:

$$f(x) = \frac{1}{m} \sum_{i=1}^m |Z_i^T x - B_i| = \frac{1}{m} \|Ax - B\|_1$$

Where rows of matrix A include Z_i 's. In this way, we can better define function f in python. Now we need to calculate sub gradient of function f. In order to calculate sub gradient of function f, we first try to gain a relation for sub gradient of norm 1 function. As we learned in the class:

$$v_i \in \partial g(x) \rightarrow v_i = \begin{cases} 1 & x_i > 0 \\ -1 & x_i < 0 \\ [-1, 1] & x_i = 0 \end{cases}$$

Now we use above formula to calculate sub gradient of our function:

$$\partial f(x) = \partial \left(\frac{1}{m} g(Ax - B) \right) = \frac{1}{m} A^T \partial g(Ax - B)$$

In fact we achieved above formula using chain rule for gradient, which is extendable to sub gradients as well. Using these relations, we implement a python code to calculate sub gradient. Finally using sub gradient we can implement Sub-Gradient Descent algorithm:

```
1 # functio below calculates subgradient of the function:
2
3 def norm1_subgrad(x):
4
5     g = np.zeros_like(x)
6     for i in range(x.size):
7         if x[i]>0:
8             g[i] = 1
9         else:
10            g[i] = -1
11
12     return g
13
14 def subgrad(x, A, B): # This function gives subgradient of f(x)
15     m = B.size
16     return 1/m * A.T @ norm1_subgrad(A @ x - B)
17
```

Now we use Sub-Gradient Descent algorithm and update values of x_t :

```
T = 4000
x0 = np.zeros(n)
eta = np.array([0.01, 0.1, 1, 10])

# Implementing Sub-Gradient Decent algorithm:

plt.figure(figsize = (18, 12))

for i in range (eta.size):
    x_new = x0
    lst = []
    for j in range(T+1):

        diff = obj_func(x_new, A, B)-obj_func(x_opt, A, B)
        lst.append(diff)
        x_new -= eta[i]*subgrad(x_new, A, B)

    plt.subplot(int(np.ceil(np.sqrt(eta.size))), int(np.ceil(np.sqrt(eta.size))), i + 1)
    plt.plot(np.arange(T + 1), lst, linewidth=2)
    plt.xlabel(r'$t$')
    plt.ylabel(r'$f(\mathbf{x}^{(t)}) - f(\mathbf{x}^{*})$')
    plt.title(r'$\eta = \{ \eta \}$.format(eta=eta[i])')
    plt.grid(True)

plt.suptitle(r'Sub-Gradient Descent On $f(\mathbf{x})$ with $\mathbf{x}_0 = \mathbf{0}$')
plt.show()
```

Now, we can see the results of updating x_t 's with different coefficients of η . Note that we used CVXPY library to calculate the minimum of the function which is denoted as $f(x^*)$ in the code. Running the code, we get the following results:

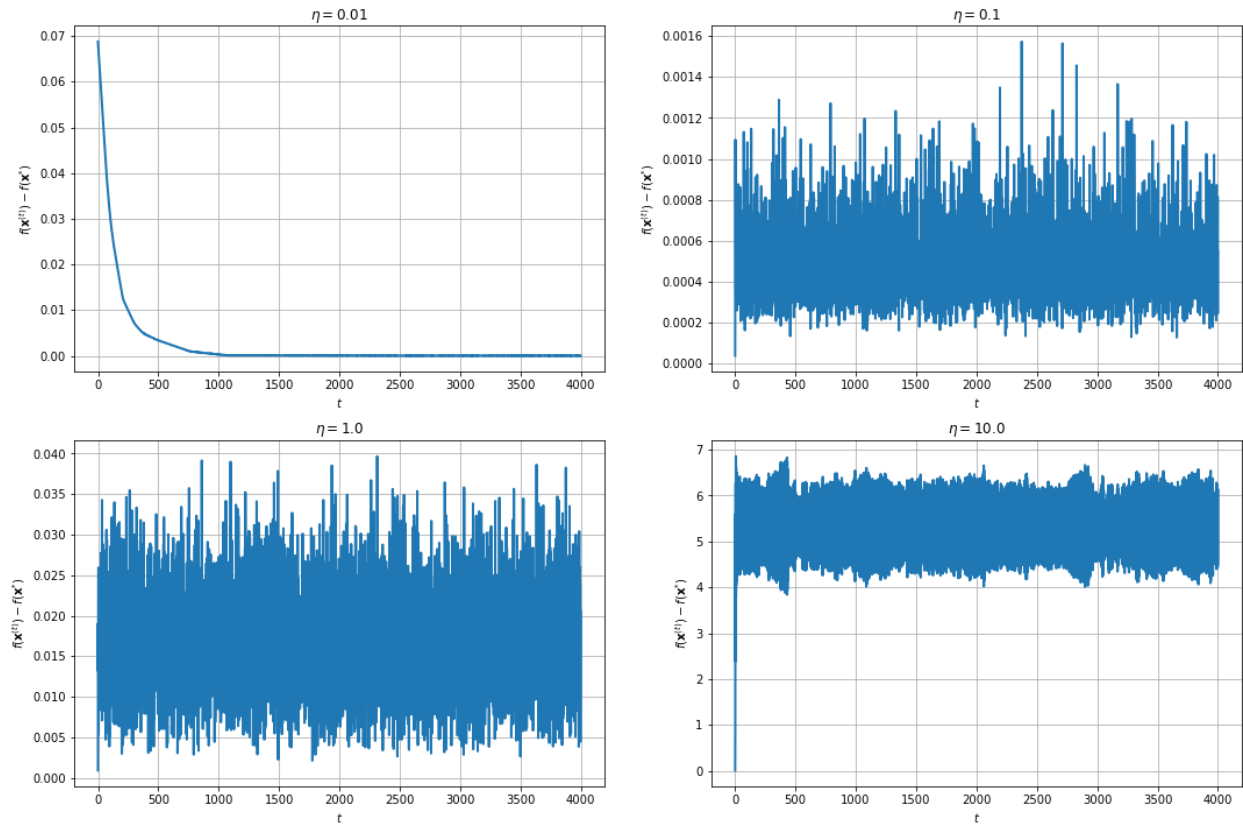
optimal X is :

```
[[-0.1667894 ]
 [-0.22008943]
 [ 0.13805849]
 [ 0.09134012]
 [-0.11563373]
 [ 0.04394651]
 [-0.15485976]
 [-0.35966015]
 [-0.04665486]
 [-0.06390327]]
```

$f(x^*) = [0.79011531]$

$f(x_0) = [0.85887246]$

Sub-Gradient Descent On $f(\mathbf{x})$ with $\mathbf{x}_0 = \mathbf{0}$



Which shows, selecting small enough value for η is essential for appropriate convergence of the Sub-Gradient Descent algorithm. As you can see $\eta = 0.01$ is the best choice for pretty good convergence in SGD algorithm.

Simulation 2:

We calculate minimum value of the function until current iteration:

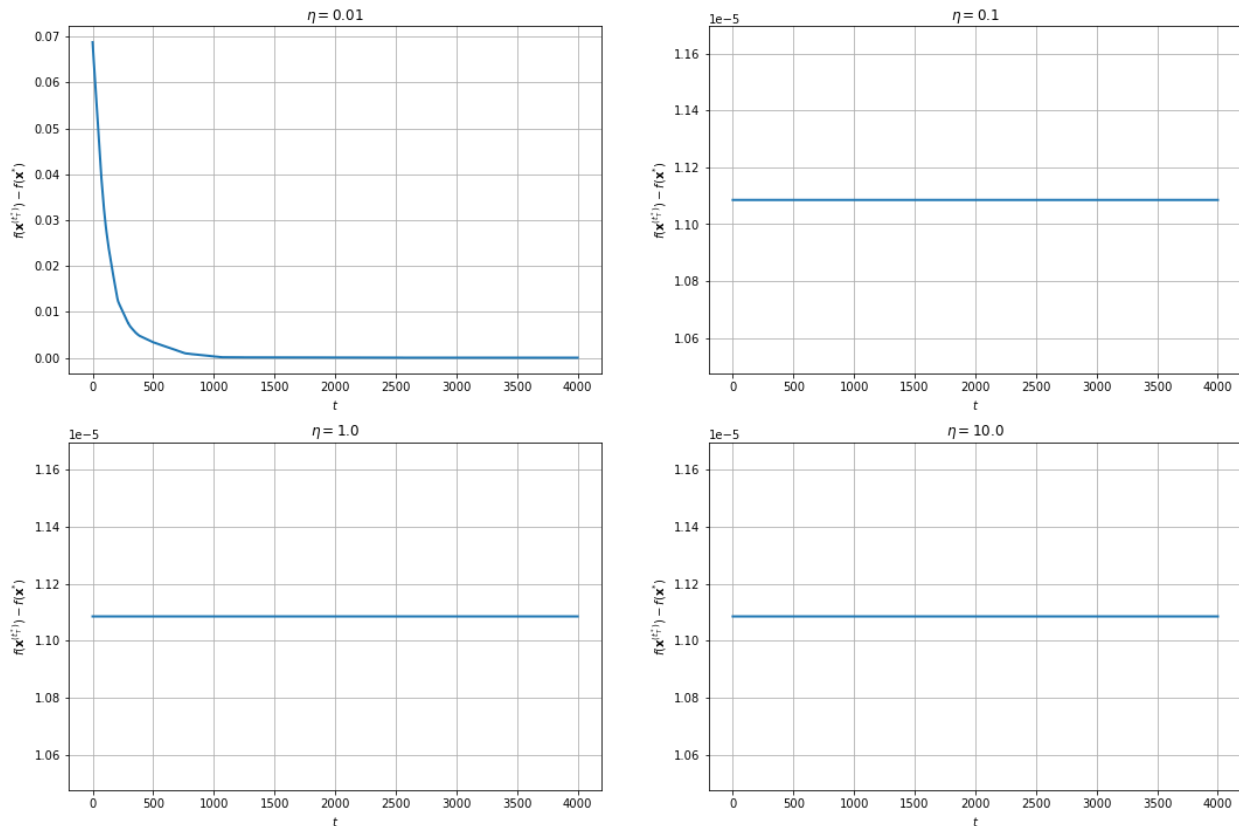
```

2 x0 = np.zeros(n)          # Initial value
3 f0 = obj_func(x0, A, B)
4
5 difference = []           # saves minimum between f(x) and f0 for every x from the begining of the algorithm until current step
6 eta = np.array([0.01, 0.1, 1, 10])
7
8 plt.figure(figsize = (18,12))
9 for i in range(eta.size):
10     x_new = x0
11     difference = []
12     for t in range(T+1):
13
14         f_xnew = obj_func(x_new, A, B)
15         f0 = np.min((f_xnew, f0))
16         x_new -= eta[i]*subgrad(x_new, A, B)
17         difference.append(f0 - obj_func(x_opt, A, B))    #f(x(t*)) - f(x*)
18
19     plt.subplot(int(np.ceil(np.sqrt(eta.size))), int(np.ceil(np.sqrt(eta.size))), i + 1)
20     plt.plot(np.arange(T + 1), difference, linewidth=2)
21     plt.xlabel(r'$t$')
22     plt.ylabel(r'$f(\mathbf{x}^{(t^*)}) - f(\mathbf{x}^{(t^*)})$')
23     plt.title(r'$\eta = \{eta\}$.format(eta=eta[i])')
24     plt.grid(True)
25
26 plt.suptitle(r'Sub-Gradient Descent On $f(\mathbf{x})$ with $\mathbf{x}_0 = \mathbf{0}$')
27 plt.show()

```

Which leads us to the following results:

Sub-Gradient Descent On $f(\mathbf{x})$ with $\mathbf{x}_0 = \mathbf{0}$



Note that for $\eta = \{0.1, 1, 10\}$ the real plot would be something like L word in English. Because it has been plotted bigger, it is not quite observable.

Simulation 3:

The code below in python implements Stochastic-Gradient Descent algorithm:

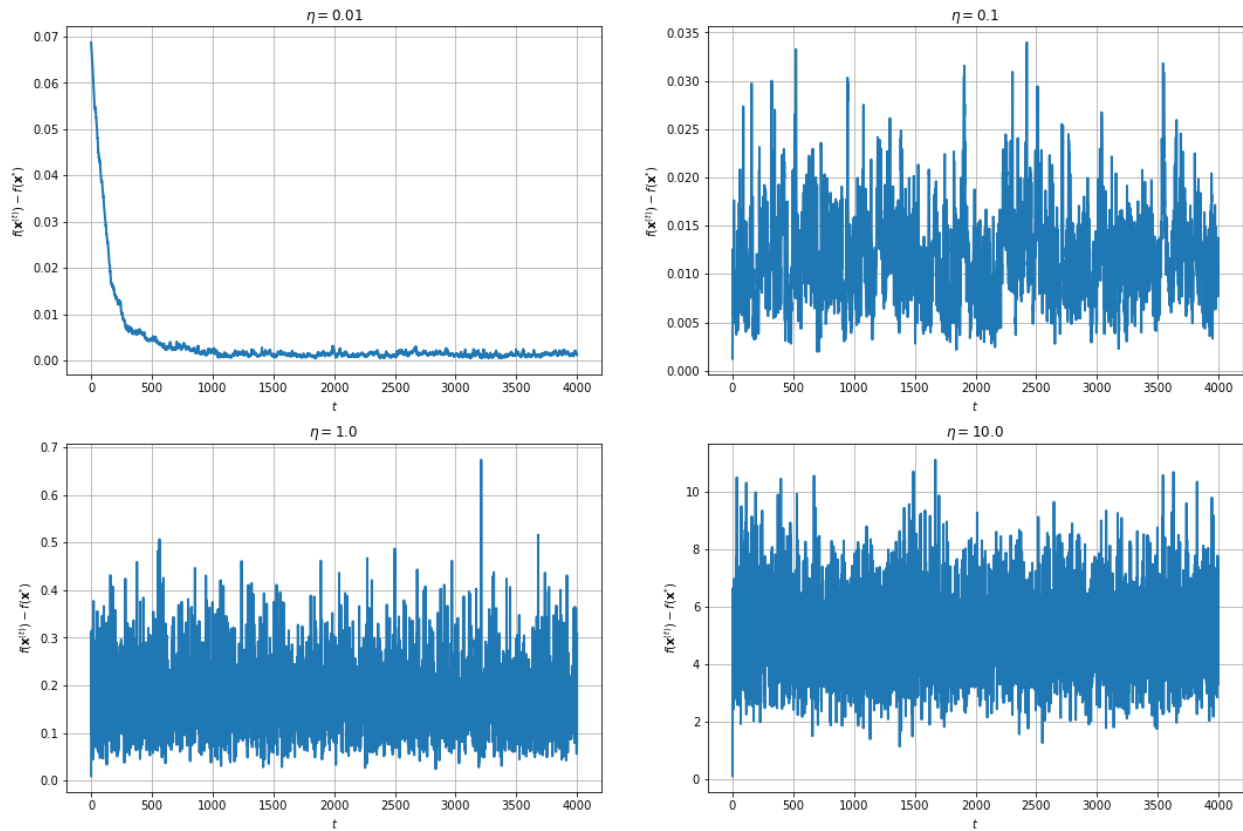
```

2 batch_size = 30
3 T = 4000
4 x0 = np.zeros(n)
5 eta = np.array([0.01, 0.1, 1, 10])
6 # Implementing Stochastic Gradient Decent algorithm:
7
8 plt.figure(figsize = (18, 12))
9
10 for i in range (eta.size):
11     x_new = x0
12     lst = []
13     for j in range(T+1):
14
15         diff = obj_func(x_new, A, B)-obj_func(x_opt, A, B)
16         lst.append(diff)
17         random_samples = np.random.choice(m, batch_size)
18
19         # choose a random sample set from dataset
20
21         A_t = A[random_samples, :]
22         B_t = B[random_samples]
23         V_t = subgrad(x_new, A_t, B_t) # random subgradient where expected value of V_T belongs to the subgradient of f(X_t)
24         x_new -= eta[i]*V_t
25
26     plt.subplot(int(np.ceil(np.sqrt(eta.size))), int(np.ceil(np.sqrt(eta.size))), i + 1)
27     plt.plot(np.arange(T + 1), lst, linewidth=2)
28     plt.xlabel(r'$t$')
29     plt.ylabel(r'$f(\mathbf{x}^{(t)}) - f(\mathbf{x}^*)$')
30     plt.title(r'$\eta = \{ \eta \}$.format(eta=eta[i])')
31     plt.grid(True)
32
33 plt.suptitle(r'Stochastic Gradient Descent On $f(\mathbf{x})$ with $\mathbf{x}_0 = \mathbf{0}$, batch size = 30$')
34 plt.show()
35

```

Now we can discuss the results:

Stochastic Gradient Descent On $f(\mathbf{x})$ with $\mathbf{x}_0 = \mathbf{0}$, batch size = 30

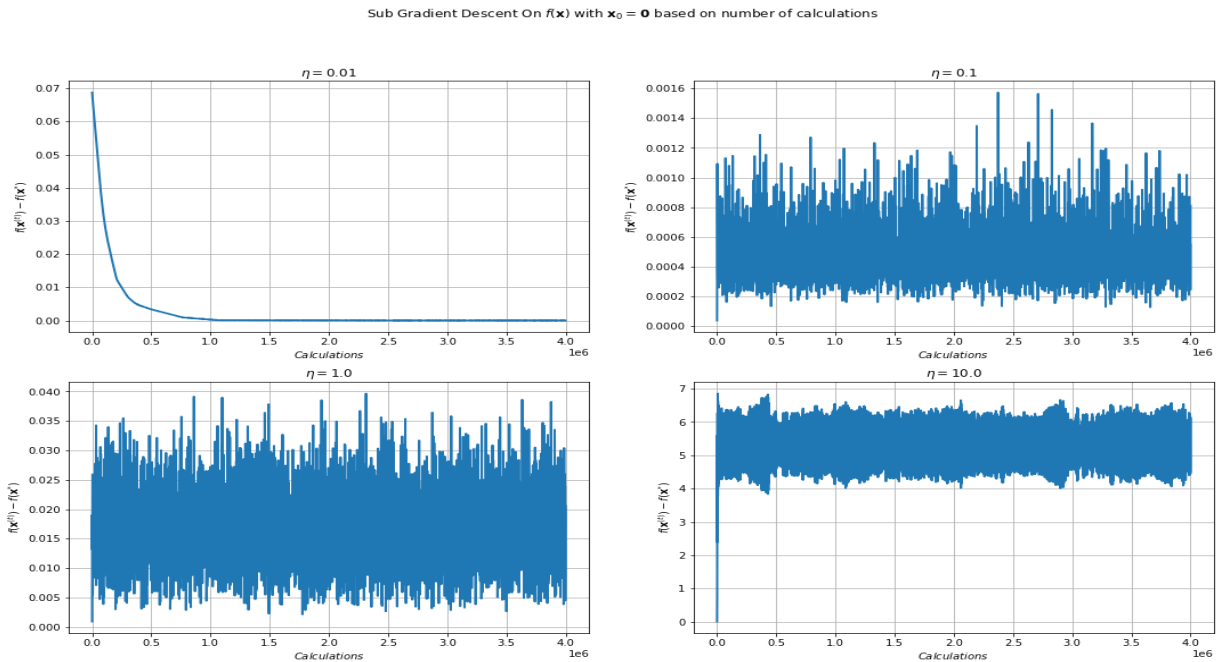


As a matter of difference, we can mention convergence speed and code's speed in running. Using Stochastic-Gradient Descent algorithm will result in faster execution of algorithm, but slightly decreases the convergence speed.

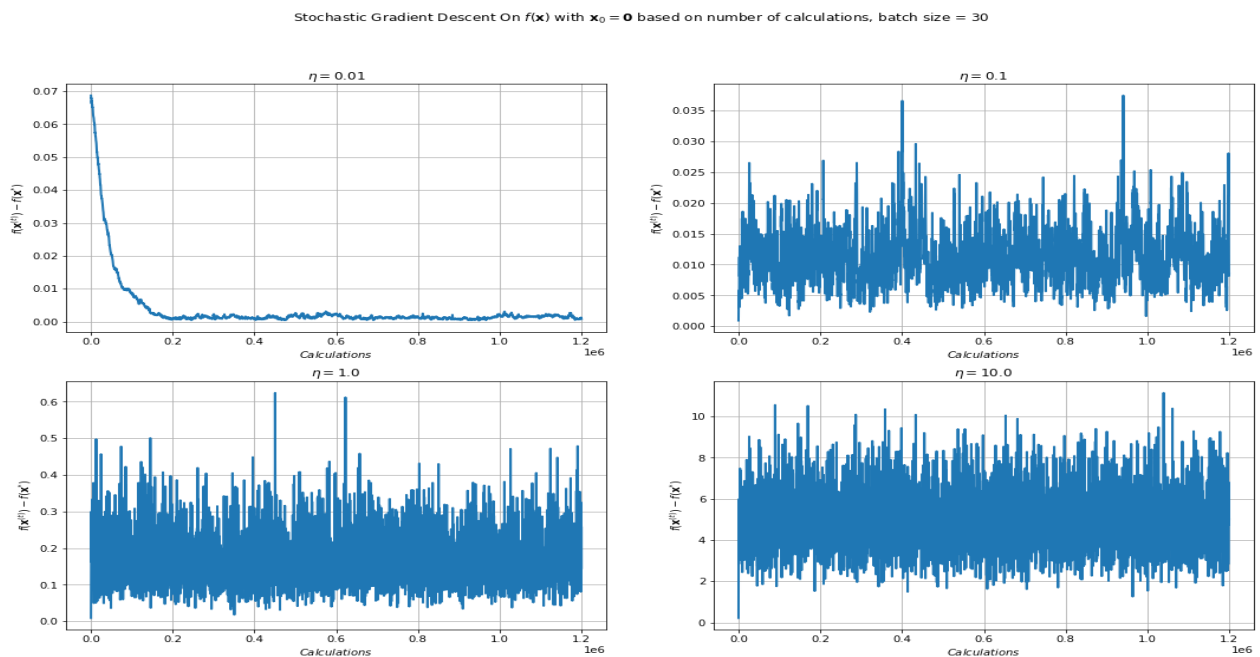
Simulation 4:

Let's see both results simultaneously.

Below we can see Sub-Gradient Descent versus number of calculations:



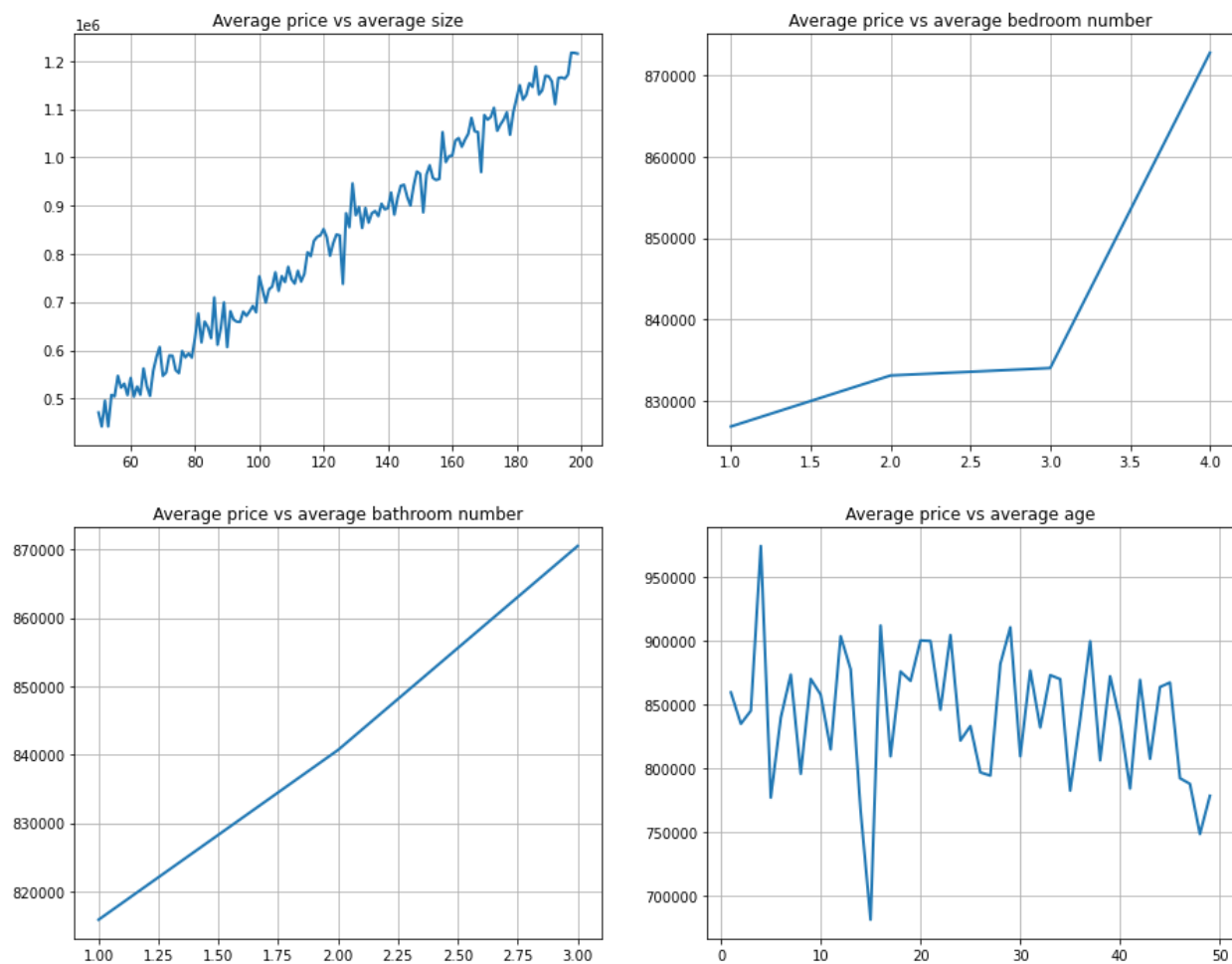
Below we can see Stochastic-Gradient Descent versus number of calculations:



It is obvious that Stochastic-Gradient Descent converges in lower number of calculations. This result totally makes sense with the previous results in simulation 3 which we mentioned that Stochastic-Gradient Descent converges to the final value in lower iterations. These results are for $\eta = 0.01$

Simulation 5:

Let's take a look at plots that demonstrate relation between a feature of a house and price of the house:



Take size of the house as a matter of illustration. As the size of the house increases, the price of the house increases proportionally. This is true for number of bathrooms as well. On the other hand, age of the house is not that important and lower age of the construction, doesn't necessarily

guarantee higher price of the house, meaning that the relation between average age of the house and the average price is not proportional. Note that repetitions of some features can be a challenge. For example there are lots of houses with 1 bedrooms which vary on price. In order to handle this issue we considered the average value of the bedrooms. That why you can see values such as 2.25 in the bathroom section.

Simulation 6:

Below code, divides data into two parts. Train data and test data. The ratio of this separation is optional and we took 0.75 of the whole data as train data and rest of the data would be test data. We did this using **sample** function in python to shuffle data and then split data into two groups.

```
1  # Splitting the data to train and test
2
3  sampled_house_data = house_data.sample(frac=1).to_numpy()      # we shuffle data using sample function
4
5  ratio = 0.75
6  total_rows = sampled_house_data.shape[0]
7  train_size = int(total_rows*ratio)
8  train_data = sampled_house_data[:train_size, :]               # This would be the training data
9  test_data = sampled_house_data[train_size:, :]                # This would be the test data
10
11
12  train_x = train_data[:, :-1]
13  train_y = train_data[:, -1].reshape(-1, 1)
14  test_x = test_data[:, :-1]
15  test_y = test_data[:, -1].reshape(-1, 1)
16
17  # Normalizing train data
18  train_x_means = train_x.mean(axis=0)
19  train_x_stds = train_x.std(axis=0)
20  train_y_means = train_y.mean(axis=0)
21  train_y_stds = train_y.std(axis=0)
22
23  train_x_normalized = (train_x - train_x_means) / train_x_stds
24  train_y_normalized = (train_y - train_y_means) / train_y_stds
25
26  test_x_normalized = (test_x - train_x_means) / train_x_stds
27  test_y_normalized = (test_y - train_y_means) / train_y_stds
28
```

Simulation 7:

Below function, calculates MSE:

```
1 # Function below calculates MSE:
2
3 def MSE_Calculator(X, y, w, b):
4     m, n = X.shape
5     one = np.ones((m, 1))
6     y_hat = X @ w + b * one
7     mse = ((y - y_hat)**2).sum()
8     return mse / m
9
```

Simulation 8:

Function below calculates SGD:

```
1
2 def SGD(X, y, w, b, lr=1e-3):
3     m, n = X.shape
4     one = np.ones((m, 1))
5     # Calculate gradients
6
7     w_grad = 2*X.T @ X @ w + 2*b * X.T @ one - 2*X.T @ y
8     b_grad = 2*one.T @ X @ w + 2*b* one.T @ one - 2*one.T @ y
9
10    w = w - lr * w_grad
11    b = b - lr * b_grad
12
13    return w, b
14
```

Simulation 9:

Take following specification:

number of epochs = 30

batch_size = 64

learning rate = 10^{-4}

```
1 2 NUM_EPOCHS = 30
3 LEARNING_RATE = 1e-4
4 BATCH_SIZE = 64
5
6 train_errors = []
7 test_errors = []
8
9 m, n = train_x.shape
10
11 # Initialize model parameters
12
13 w = np.random.rand(n, 1)
14 b = np.random.rand(1)
15
16 for epoch in range(1, NUM_EPOCHS+1):
17     w, b = iteration(
18         X=train_x_normalized,
19         y=train_y_normalized,
20         w=w,
21         b=b,
22         batch_size=BATCH_SIZE,
23         lr=LEARNING_RATE
24     )
25     mse_train = MSE_Calculator(train_x_normalized, train_y_normalized, w, b)
26     mse_test = MSE_Calculator(test_x_normalized, test_y_normalized, w, b)
27     train_errors.append(mse_train)
28     test_errors.append(mse_test)
29     print(f'epoch {epoch}, train mse: {np.round(mse_train, 3)}, test mse: {np.round(mse_test, 3)}')
30     print()
31
```

```
1
2 fig = plt.figure(figsize=(8, 6))
3 ax = fig.add_subplot(1, 1, 1)
4 ax.plot(range(1, NUM_EPOCHS+1), train_errors, label='train error')
5 ax.plot(range(1, NUM_EPOCHS+1), test_errors, label='test error')
6 ax.legend();
7 ax.grid(True)
8 ax.set_title('Test and train MSEs after each epoch')
9 ax.set_xlabel('Epochs')
10 ax.set_ylabel('MSE')
```


And the result of the code would be:

