

Theory Questions

1.

a. Let there be f, A and b as required.

Let there be $x, y \in \mathbb{R}^n, \lambda \in [0, 1]$.

$$\begin{aligned} g(\lambda x + (1 - \lambda)y) &= f(A\lambda x + A(1 - \lambda)y + b) = \\ f(\lambda Ax + \lambda b + (1 - \lambda)Ay + (1 - \lambda)b) &= f(\lambda(Ax + b) + (1 - \lambda)(Ay + b)) \leq \\ \lambda f(Ax + b) + (1 - \lambda)f(Ay + b) &= \lambda g(x) + (1 - \lambda)g(y) \end{aligned}$$

b. Let there be $\{f_i\}_{i=1}^m$ as required.

Let there be $x, y \in \mathbb{R}^n, \lambda \in [0, 1]$.

$$\begin{aligned} g(\lambda x + (1 - \lambda)y) &= \max_i f_i(\lambda x + (1 - \lambda)y) \\ &\leq \max_i \lambda f_i(x) + (1 - \lambda)f_i(y) \\ &\leq \max_i \lambda f_i(x) + \max_i (1 - \lambda)f_i(y) = \lambda g(x) + (1 - \lambda)g(y) \end{aligned}$$

c. We will show that ℓ_{\log} is convex:

$$\frac{d\ell_{\log}}{dz} = \frac{-e^{-z}}{(1 + e^{-z})\ln 2}$$

$$\begin{aligned} \frac{d^2\ell_{\log}}{dz^2} &= \frac{e^{-z}(1 + e^{-z})\ln 2 - (-e^{-z})\ln 2(-e^{-z})}{((1 + e^{-z})\ln 2)^2} \\ &= \frac{\ln 2e^{-z} + \ln 2e^{-2z} - \ln 2e^{-2z}}{((1 + e^{-z})\ln 2)^2} \\ &= \frac{\ln 2e^{-z}}{((1 + e^{-z})\ln 2)^2} \end{aligned}$$

Because $\ln 2 > 0$, and for every $z \in \mathbb{R}, e^{-z} > 0$, the second derivative $\frac{d^2\ell_{\log}}{dz^2}$ is positive for every $z \in \mathbb{R}$, hence ℓ_{\log} is convex.

In the first article of this question, we showed that for any convex function $f: \mathbb{R}^n \rightarrow \mathbb{R}, A \in \mathbb{R}^{n \times n}$, and $b \in \mathbb{R}^n, g(x) = f(Ax + b)$ is convex.

Therefore, if we take $y \in \mathbb{R}, x \in \mathbb{R}^n$, and set b as the zero vector in \mathbb{R}^n we can conclude that $f: \mathbb{R}^n \rightarrow \mathbb{R}, f(w) = \ell_{\log}(yw \cdot x)$ is convex.

Because this is a composition of convex function with affine function, as in article (a).

2. We will prove that for $w_{hinge}^* \in \operatorname{argmin}_{w \in \mathbb{R}^d} \sum_{i=1}^n \ell_{hinge}(y_i w \cdot x_i)$ satisfies that $\forall i \operatorname{sign}(w_{hinge}^* \cdot x_i) = y_i$.

First, we will prove the following **proposition**:

for any $r \in \mathbb{R}$, $\ell_{zo}(r) \leq \ell_{hinge}(r)$, when $\ell_{zo}(r) = \begin{cases} 0 & r > 0, \\ 1 & r \leq 0 \end{cases}$,
 $\ell_{hinge}(r) = \max\{0, 1 - r\}$.

Prove:

If $r > 0$, then $\ell_{hinge}(r) = \max\{0, 1 - r\} \geq 0 = \ell_{zo}(r)$

If $r \leq 0$, then $\ell_{hinge}(r) = \max\{0, 1 - r\} \geq 1 = \ell_{zo}(r)$

Therefore, for any $r \in \mathbb{R}$, $\ell_{zo}(r) \leq \ell_{hinge}(r)$.

Now we focus on proving that w_{hinge}^* is also optimal for ℓ_{zo} .

Let there be $i \in \{1, \dots, n\}$, $c \in \mathbb{R}$.

Let us see how cw^* behaves: $y_i cw^* \cdot x_i = cy_i w^* \cdot x_i > 0$

Because $y_i w^* \cdot x_i > 0$ for all i .

If we take $c > \frac{1}{y_i w^* \cdot x_i}$, we notice that $y_i cw^* \cdot x_i = cy_i w^* \cdot x_i > \frac{1}{y_i w^* \cdot x_i} \cdot y_i w^* \cdot x_i = 1$,

Hence $\ell_{hinge}(y_i cw^* \cdot x_i) = 0$

Using the hint, if we take $c \rightarrow \infty$, we get this constraint ($c > \frac{1}{y_i w^* \cdot x_i}$) for all i .

So, we conclude that for large enough c , $\forall i \ell_{hinge}(y_i cw^* \cdot x_i) = 0$,

So, $\sum_{i=0}^n \ell_{hinge}(y_i cw^* \cdot x_i) = 0$.

Therefore,

$$\begin{aligned} 0 &\leq \sum_{i=0}^n \ell_{hinge}(y_i w_{hinge}^* \cdot x_i) \leq 0 \\ &\Rightarrow \sum_{i=0}^n \ell_{hinge}(y_i w_{hinge}^* \cdot x_i) = 0 \\ &\Rightarrow \forall i \ell_{hinge}(y_i w_{hinge}^* \cdot x_i) = 0 \end{aligned}$$

Where the last statements follow from the fact that ℓ_{hinge} is non-negative.

Because the hinge lost bounds the zero-one loss, we conclude that

$$\forall i \ell_{zo}(y_i w_{hinge}^* \cdot x_i) = 0$$

And this happens only if

$$\forall i \operatorname{sign}(w_{hinge}^* \cdot x_i) = y_i$$

3.

a.

Let there be $y \in \mathbb{R}^d$, $x = \Pi_K(y)$ it's projection on K .

Because K is convex, for any $\lambda \in (0, 1)$, $z \in K$ we have that

$$\lambda z + (1 - \lambda)x = x + \lambda(z - x) \in K$$

Therefore, because $\|x - y\|_2 = \min_{z \in K} \|y - z\|$

$$\begin{aligned} \|x - y\|_2^2 &\leq \|x + \lambda(z - x) - y\|_2^2 \\ &= \|x - y\|_2^2 + 2\lambda\langle x - y, z - x \rangle + \lambda^2\|z - x\|_2^2 \end{aligned}$$

So, we conclude that $2\lambda\langle x - y, z - x \rangle \geq -\lambda^2\|z - x\|_2^2$.

From taking $\lambda \rightarrow 0$ we can see that $\langle x - y, z - x \rangle \geq 0$

So,

$$\begin{aligned} \|y - z\|_2^2 &= \|z - y\|_2^2 = \|z - x + x - y\|_2^2 \\ &= \|z - x\|_2^2 + 2\langle x - y, z - x \rangle + \|x - y\|_2^2 \\ &\geq \|z - x\|_2^2 \end{aligned}$$

Where in the last inequality we used that $\langle x - y, z - x \rangle \geq 0$, $\|x - y\|_2^2 \geq 0$.

and this yields us that $\|y - z\|_2 \geq \|x - z\|_2$

b.

Let there be $x^* \in \operatorname{argmin}_{x \in K} f(x)$, $\eta_t = \frac{\epsilon}{G^2}$ and $T = \frac{B^2 G^2}{\epsilon^2}$ as in the Theorem.

I assume that $\|x^*\|_2 \leq B$, $\|\nabla f(x_t)\|_2 \leq G$ for every $x \in K$, and $x_1 = \mathbf{0}$ (like what was presented in Lecture #5).

Using Jensen inequality, we obtain that

$$\begin{aligned} f(\bar{x}) - f(x^*) &= f\left(\sum_{t=1}^T x_t\right) - f(x^*) \leq \frac{1}{T} \sum_{t=1}^T f(x_t) - f(x^*) \\ &= \frac{1}{T} \sum_{t=1}^T (f(x_t) - f(x^*)) \end{aligned}$$

Because f is convex, we get that for any t ,

$$f(x_t) - f(x^*) \leq \langle \nabla f(x_t), x_t - x^* \rangle$$

Thus,

$$f(\bar{x}) - f(x^*) \leq \frac{1}{T} \sum_{t=1}^T \langle \nabla f(x_t), x_t - x^* \rangle$$

And because $y_{t+1} = x_t - \eta_t \nabla f(x_t)$,
we get that

$$f(\bar{x}) - f(x^*) \leq \frac{1}{T} \sum_{t=1}^T \langle \nabla f(x_t), x_t - x^* \rangle \leq \frac{1}{T} \sum_{t=1}^T \frac{1}{\eta_t} \langle x_t - y_{t+1}, x_t - x^* \rangle$$

Using completing the square we get that

$$(a - b)(a - c) = \frac{(a - b)^2 + (a - c)^2 - (b - c)^2}{2}$$

Therefore,

$$\begin{aligned} & \frac{1}{\eta_t} \langle x_t - y_{t+1}, x_t - x^* \rangle \\ &= \frac{1}{2\eta_t} (\|x_t - y_{t+1}\|_2^2 + \|x_t - x^*\|_2^2 - \|y_{t+1} - x^*\|_2^2) \end{aligned}$$

And if we again use that $y_{t+1} = x_t - \eta_t \nabla f(x_t)$, we get

$$= \frac{1}{2\eta_t} (\|x_t - x^*\|_2^2 - \|y_{t+1} - x^*\|_2^2) + \frac{\eta_t}{2} \|\nabla f(x_t)\|_2^2$$

From article (a), we obtain that

$$\|x_{t+1} - x^*\|_2 \leq \|y_{t+1} - x^*\|_2$$

And this is the same as

$$-\|x_{t+1} - x^*\|_2 \geq -\|y_{t+1} - x^*\|_2$$

So, if we plug this into the upper formula, we get that

$$\begin{aligned} & \frac{1}{2\eta_t} (\|x_t - x^*\|_2^2 - \|y_{t+1} - x^*\|_2^2) + \frac{\eta_t}{2} \|\nabla f(x_t)\|_2^2 \\ & \leq \frac{1}{2\eta_t} (\|x_t - x^*\|_2^2 - \|x_{t+1} - x^*\|_2^2) + \frac{\eta_t}{2} \|\nabla f(x_t)\|_2^2 \end{aligned}$$

For conclusion, we got that

$$\frac{1}{\eta_t} \langle x_t - y_{t+1}, x_t - x^* \rangle \leq \frac{1}{2\eta_t} (\|x_t - x^*\|_2^2 - \|x_{t+1} - x^*\|_2^2) + \frac{\eta_t}{2} \|\nabla f(x_t)\|_2^2$$

So,

$$\begin{aligned} f(\bar{x}) - f(x^*) & \leq \frac{1}{T} \sum_{t=1}^T \frac{1}{\eta_t} \langle x_t - y_{t+1}, x_t - x^* \rangle \\ & \leq \frac{1}{T} \sum_{t=1}^T \left(\frac{1}{2\eta_t} (\|x_t - x^*\|_2^2 - \|x_{t+1} - x^*\|_2^2) + \frac{\eta_t}{2} \|\nabla f(x_t)\|_2^2 \right) \\ & = \frac{1}{2T\eta_t} (\|x_1 - x^*\|_2^2 - \|x_{T+1} - x^*\|_2^2) + \frac{\eta_t}{2T} \|\nabla f(x_t)\|_2^2 \\ & \leq \frac{1}{2T\eta_t} \|x_1 - x^*\|_2^2 + \frac{\eta_t}{2T} \|\nabla f(x_t)\|_2^2 \leq \frac{1}{2T\eta_t} \|x^*\|_2^2 + \frac{\eta_t}{2} \|\nabla f(x_t)\|_2^2 \\ & \leq \frac{B^2}{2T\eta_t} + \frac{\eta_t}{2} G^2 = \frac{B^2}{2} \frac{G^2 \epsilon^2}{\epsilon B^2 G^2} + \frac{\epsilon}{2G^2} G^2 = \frac{\epsilon}{2} + \frac{\epsilon}{2} = \epsilon \end{aligned}$$

4. Let there be $f: \mathbb{R}^n \rightarrow \mathbb{R}$ as in the question (differential, non-negative and β smooth).
Let us consider x_t, x_{t+1} points in GD, where $\eta < \frac{2}{\beta}$.

Because f is β smooth:

$$\begin{aligned} f(x_{t+1}) &\leq f(x_t) + \nabla f(x_t)^T (x_{t+1} - x_t) + \frac{\beta}{2} \|x_{t+1} - x_t\|^2 \\ &= f(x_t) - \eta \nabla f(x_t)^T \nabla f(x_t) + \frac{\beta \eta^2}{2} \|\nabla f(x_t)\|^2 \\ &= f(x_t) - \|\nabla f(x_t)\|^2 \left(\eta - \frac{\beta}{2} \eta^2 \right) \end{aligned}$$

Therefore,

$$\|\nabla f(x_t)\|^2 \leq \left(\eta - \frac{\beta}{2} \eta^2 \right)^{-1} (f(x_t) - f(x_{t+1}))$$

Let's denote $S_n = \sum_{t=1}^n \|\nabla f(x_t)\|^2$ the partial sums sequence.

We obtain from the above that:

$$S_n \leq \left(\eta - \frac{\beta}{2} \eta^2 \right)^{-1} (f(x_1) - f(x_{n+1}))$$

Let us calculate when $\eta - \frac{\beta}{2} \eta^2 > 0$:

$$\eta - \frac{\beta}{2} \eta^2 > 0 \Leftrightarrow \eta \left(1 - \frac{\beta}{2} \eta \right) > 0 \Leftrightarrow 0 < \eta < \frac{2}{\beta}$$

And from our assumptions, this holds.

So, we can use that $\forall x \in \mathbb{R}^n, f(x) \geq 0$ to obtain:

$$S_n \leq \left(\eta - \frac{\beta}{2} \eta^2 \right)^{-1} f(x_1)$$

This means that S_n is bounded from above by non-negative number:

$$\left(\eta - \frac{\beta}{2} \eta^2 \right)^{-1} f(x_1)$$

We notice that $S_{n+1} \geq S_n$, because

$$S_{n+1} = S_n + \|\nabla f(x_{n+1})\|^2 \geq S_n \quad (\forall t \in \mathbb{N} \quad \|\nabla f(x_{n+1})\|^2 \geq 0)$$

This means that S_n converges because non-decreasing and bounded from above sequence converges to its supremum.

Therefore $\lim_{t \rightarrow \infty} \|\nabla f(x_t)\|^2 = 0$ from the hint, and therefore $\lim_{t \rightarrow \infty} \|\nabla f(x_t)\| = 0$

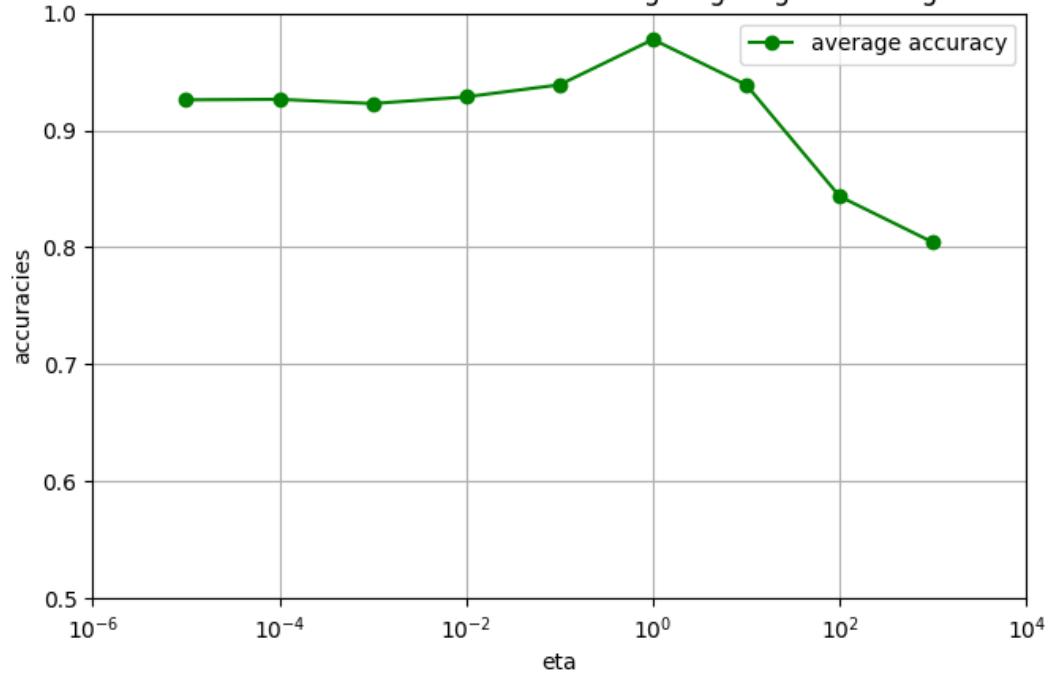
(If $\lim_{n \rightarrow \infty} a_n^2 = 0$, then $\lim_{n \rightarrow \infty} a_n = 0$).

Programming Questions

1.

a.

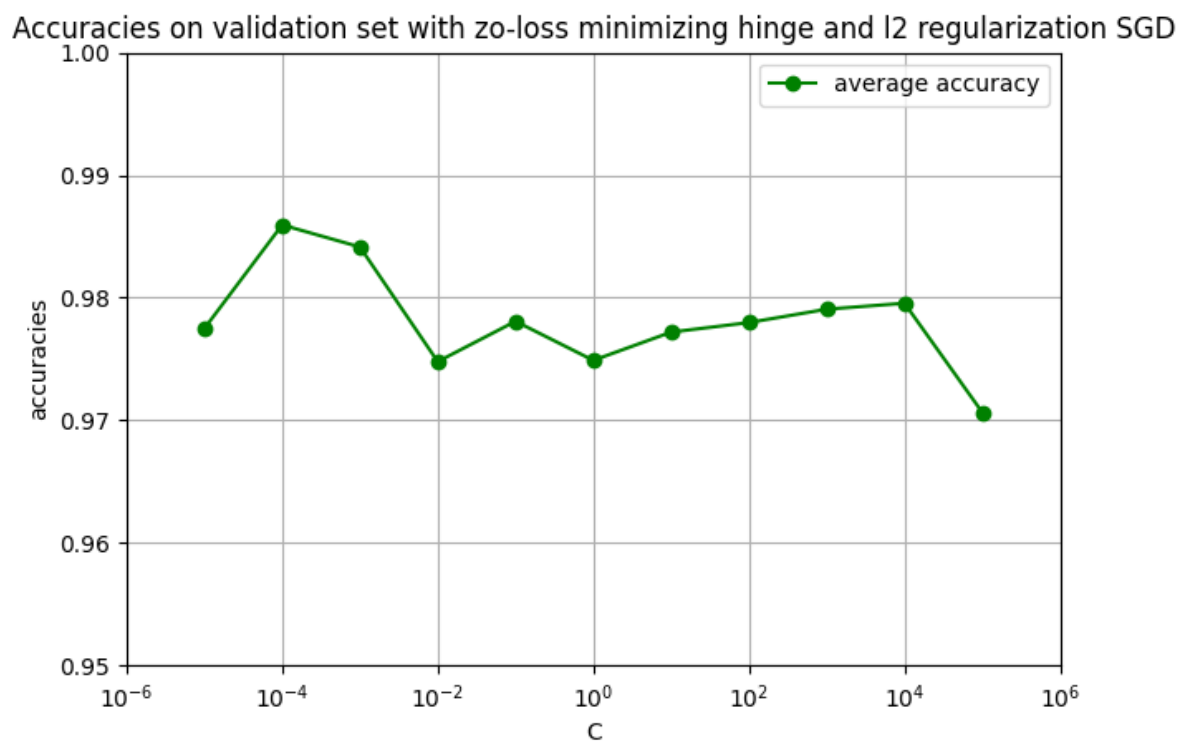
Accuracies on validation set with zo-loss minimizing hinge log with l2 regularization SGD



When η was larger than 1000 there overflow encountered in multiplication, so we tested up to $\eta = 1000$.

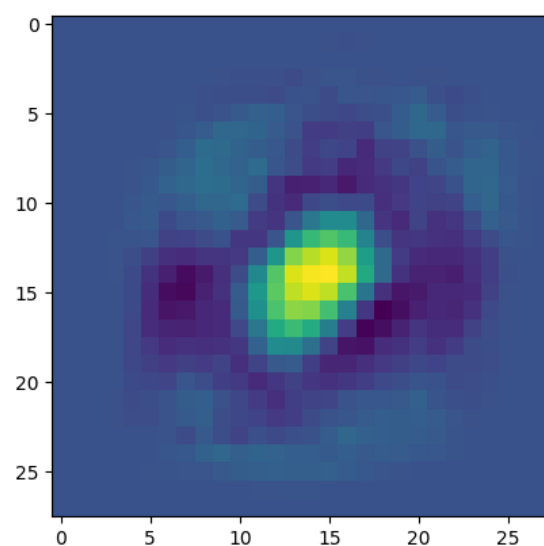
The best η is 1.

b.



The best C is 0.0001.

c.

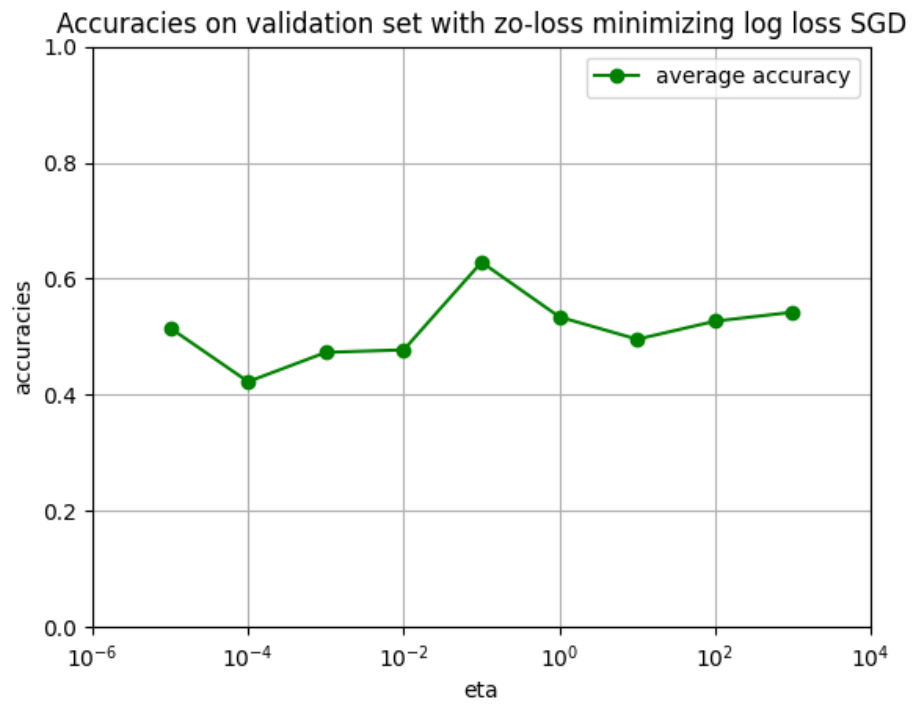


We can interpretate that the “warmer” the color, the result of SGD is denser, and when the color is “cooler”, it’s sparser.

d. The accuracy on the test dataset is 99.28352098259979%.

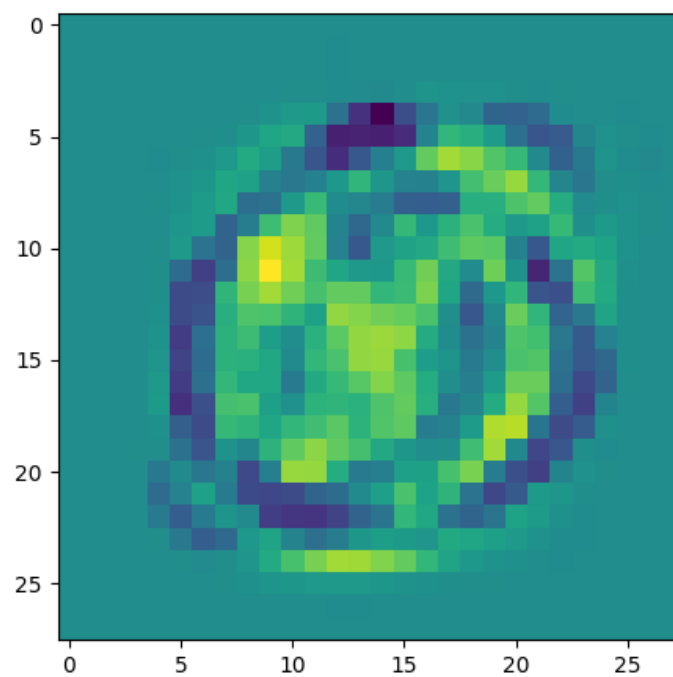
2.

a.



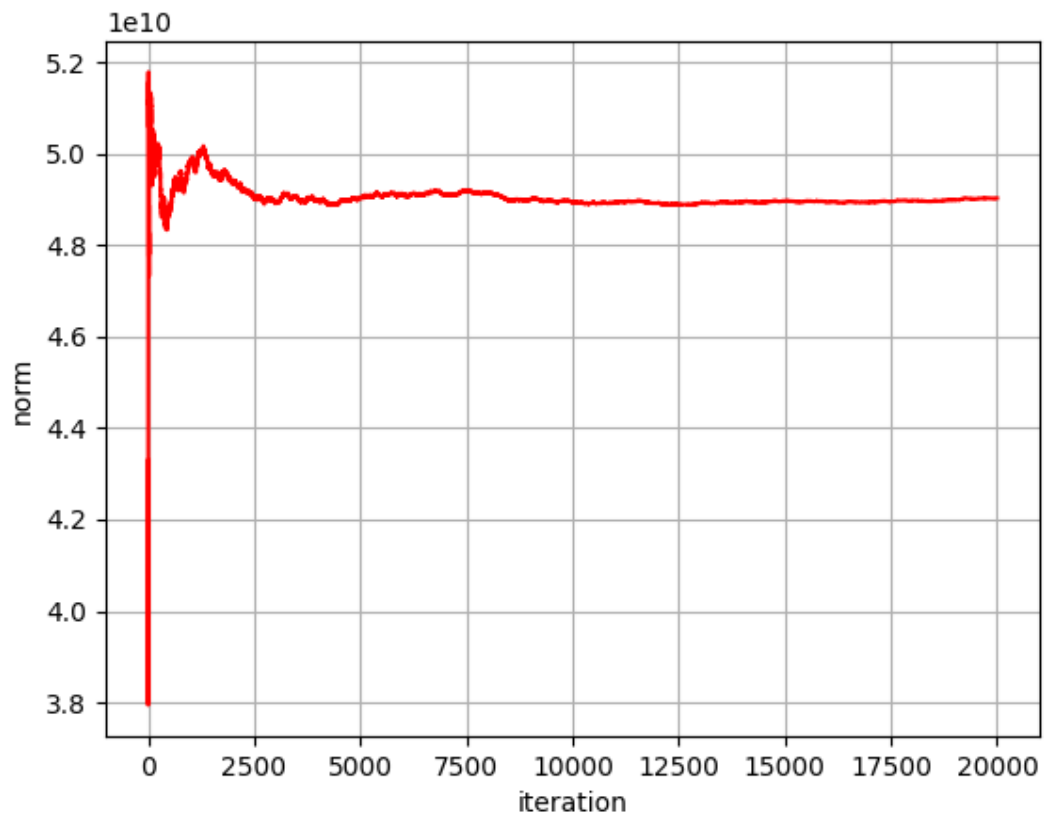
The best η is 0.1.

b.



The accuracy is 75.84442169907881%.

c.



We can see that the norms sequence converges.

We know that in expectation SGD converges to the optimal point.

Probably, this number of iterations was large enough for SGD to convert for optimum, and because of that the norms sequence convergent to the norm of the optimum.