

חלק 1:

1. מצא מטריצה $X \in \mathbb{R}^{m \times n+1}$ כך שהבעיות שקולות:

פתרון:

$$X = \begin{pmatrix} x_1^0 & x_1^1 & \dots & x_1^n \\ x_2^0 & x_2^1 & \dots & x_2^n \\ \vdots & \vdots & \ddots & \vdots \\ x_m^0 & x_m^1 & \dots & x_m^n \end{pmatrix} \text{ נגדיר}$$

X כלומר האיבר ה- i, j הוא הדגימה ה- i בחזקת j .

$$Xa = \begin{pmatrix} x_1^0 & x_1^1 & \dots & x_1^n \\ x_2^0 & x_2^1 & \dots & x_2^n \\ \vdots & \vdots & \ddots & \vdots \\ x_m^0 & x_m^1 & \dots & x_m^n \end{pmatrix} \begin{pmatrix} a_0 \\ \vdots \\ a_n \end{pmatrix} = \begin{pmatrix} x_1^0 \cdot a_0 + x_1^1 \cdot a_1 + \dots + x_1^n \cdot a_n \\ \vdots \\ x_m^0 \cdot a_0 + x_m^1 \cdot a_1 + \dots + x_m^n \cdot a_n \end{pmatrix} = \begin{pmatrix} f(x_1) \\ \vdots \\ f(x_m) \end{pmatrix}$$

נשים לב כי:

$$\Rightarrow (Xa)_i = f(x_i)$$

$$\text{לכן: } h(a) = \frac{1}{2m} \|y - Xa\|_{L2 \text{ Norm}}^2 = \frac{1}{2m} \sum_{i=1}^m (y_i - (Xa)_i)^2 = \frac{1}{2m} \sum_{i=1}^m (y_i - f(x_i))^2$$

שקולות עבור המטריצה הזאת.

$$2. \text{ נראה בעת כי } (P) = \min_{a \in \mathbb{R}^{n+1}} \left\{ h(a) = \frac{1}{2m} \sum_{i=1}^m (y_i - (Xa)_i)^2 = \frac{1}{2m} \sum_{i=1}^m (y_i - f(x_i))^2 \right\} \text{ קמורה,}$$

ונמצא פתרון סגור.

פתרון:

$$\begin{aligned} h(a) &= \frac{1}{2m} \|y - Xa\|_2^2 = (y - Xa)^T (y - Xa) = \\ &= \frac{1}{2m} (y^T - a^T X^T) (y - Xa) = \frac{1}{2m} \left(\underset{\text{scalar}}{y^T y} - \underset{\text{scalar}}{a^T X^T y} - \underset{\text{scalar}}{y^T Xa} + a^T X^T Xa \right) = \frac{1}{2m} \left(\|y\|_2^2 - (a^T X^T y) - (a^T X^T y)^T + a^T X^T Xa \right) = \\ &= \frac{1}{2m} \left(\|y\|_2^2 - 2(a^T X^T y) + a^T X^T Xa \right) \\ \nabla h(a) &= \frac{1}{2m} (-2X^T y + 2X^T Xa) \\ \nabla^2 h(a) &= \frac{1}{2m} (2X^T X) = \frac{1}{m} (X^T X) \end{aligned}$$

$X^T X$ is PSD (by Homework 1) and $m > 0 \Rightarrow \nabla^2 h(a)$ is PSD

והפתרון הסגור:

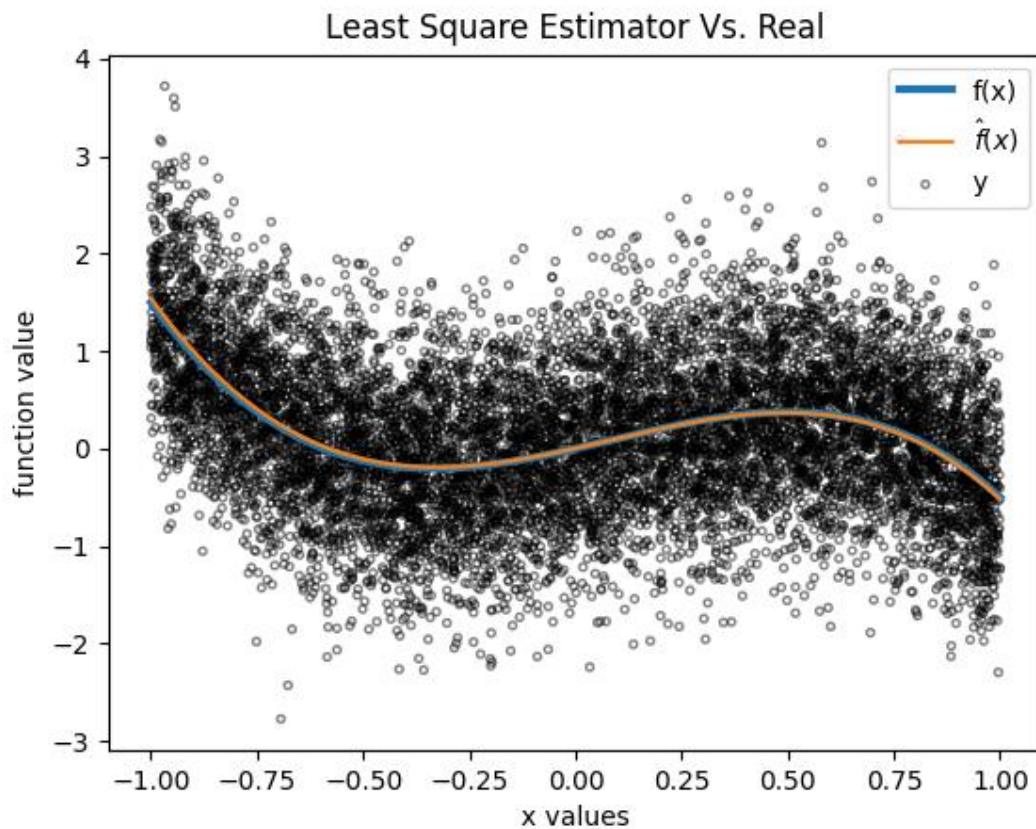
$$\begin{aligned} \nabla h(a) &= \frac{1}{2m} (-2X^T y + 2X^T Xa) = 0 \\ \Leftrightarrow X^T Xa &= X^T y \\ \Leftrightarrow a^* &= (X^T X)^{-1} X^T y \end{aligned}$$

נשים לב כי הפתרון קיים אם $X^T X$ הפיכה.

÷ 0	÷ 1	÷ 2	÷ 3
0.00294	0.97805	0.50909	-2.00785

3. המקדמים שהפונקציה שלנו החזירה הם:

4. Least Square Estimator Vs. Real



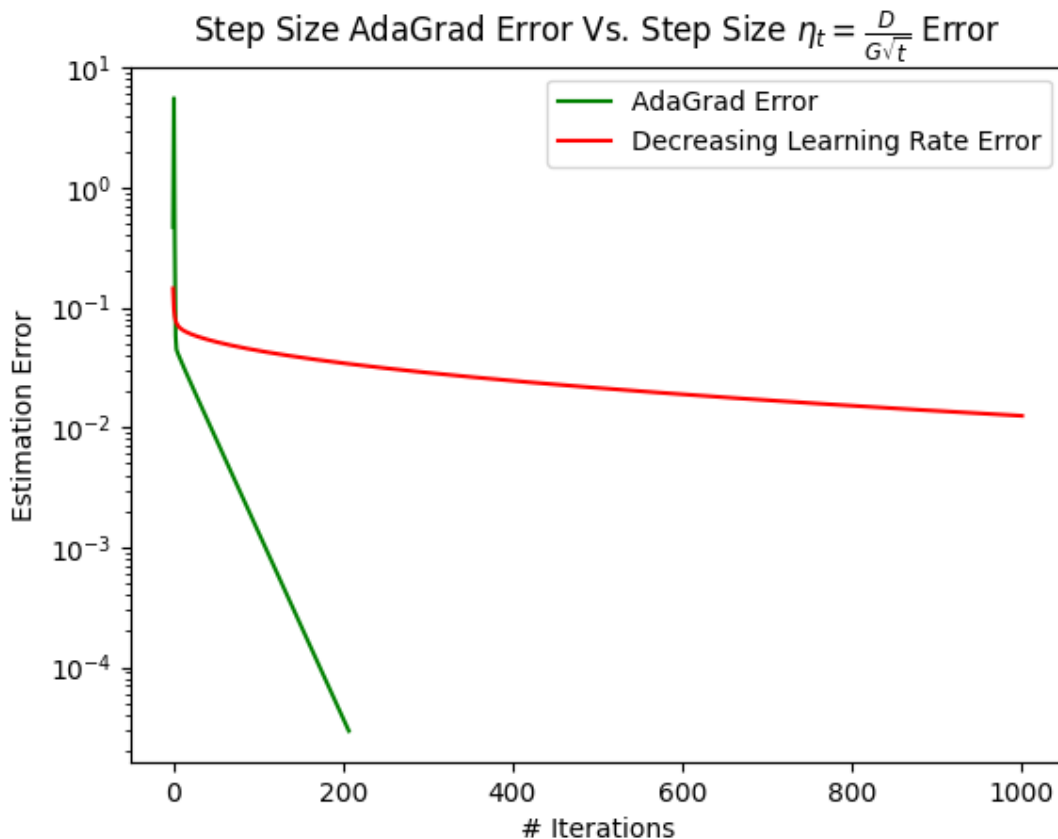
ניתן לראות מהגרף, כי הדגימות הרועשות y מופלגות בצורה גאוסית סביב הפונקציה המקורית, כלומר הסבירות לדגימה רועשת ממש (רחוקה מאוד מהפונקציה) קטנה (יש מעט דגימות שבאלה). בנוסף ניתן לראות שלמרות שהדגימות שלנו רועשות אכן שערך שמצאנו משערך היטב את הפונקציה.

חלק 2 Projected Gradient Descent

5. מהחלק הקודם של השאלה ראינו כי:

$$h(a) = \frac{1}{2m} \|y - Xa\|_2^2; \nabla h(a) = \frac{1}{2m} (-2X^T y + 2X^T Xa)$$

6. השוואה בין גודל הצעד כאשר הוא AdaGrad או $\eta_t = \frac{D}{G\sqrt{t}}$:



הרצת אלגוריתם *Projected gradient descent* המבוקש עם פרמטר סף של $\varepsilon = 0.001$.

כפי שניתן לראות מהגרף, עבור גודל צעד דועך האלגוריתם עדיין לא התכנס לטווח השגיאה הרצוי כי ראינו שלצורך כך דרושות $O\left(\frac{1}{\varepsilon^2}\right)$ איטרציות, כלומר סדר גודל של 1,000,000 איטרציות. לעומת זאת,

כאשר גודל הצעד הוא AdaGrad הביצועים הם הרבה יותר טובים, אנו משיגים שגיאה קטנה מהסף בהרבה פחות איטרציות. לפי מה שראינו בתרגול 6, AdaGrad יכול להשיג ביצועים הרבה יותר טובים כאשר הנורמה של הגראדינטים קטנה משמעותית מ-G.

7. צריך להוכיח כי קבוע החלקות של פונקציית המטרה $h(a)$ הוא: $L = \lambda_{\max}\left(\frac{X^T X}{m}\right)$.

הוכחה: נראה כי הפונקציה $h(a)$ היא בעלת גראדינט ליפשיץ עם קבוע $L = \lambda_{\max}\left(\frac{X^T X}{m}\right)$.

$$a_1, a_2 \in C$$

$$\left\|\nabla h(a_1)-\nabla h(a_2)\right\|_2^2=\frac{1}{2 m^2}\left\|-2 X^T y+2 X^T X a_1+2 X^T y-2 X^T X a_2\right\|_2^2=\frac{1}{2 m^2}\left\|2 X^T X a_1-2 X^T X a_2\right\|_2^2=\left\|\frac{X^T X}{m^2} a_1-\frac{X^T X}{m^2} a_2\right\|_2^2$$

$$let\ Z=\frac{X^TX}{m^2};Z\in R^{n+1\times n+1}$$

$$\left\|Za_1-Za_2\right\|_2^2=\left\|Z\left(a_1-a_2\right)\right\|_2^2=\left(Z\left(a_1-a_2\right)\right)^TZ\left(a_1-a_2\right)=\left(a_1-a_2\right)^TZ^TZ\left(a_1-a_2\right)\leq\lambda_{\max }\left(Z^TZ\right)\left\|a_1-a_2\right\|_2^2$$

$$, L = \lambda_{\max }\left(\frac{X^TX}{m}\right) \text{ נסמן ב-} L \geq 0 \text{ כי זאת מטריצה PSD ובפרט, כמו כן}$$

$$\lambda_{\max }\left(Z^TZ\right)=\lambda_{\max }\left(\frac{X^TX}{m} \times \frac{X^TX}{m}\right)=L \cdot L=L^2$$

$$L-\text{זאת מביוון ש-} L \text{ הוא ערך עצמי של } \frac{X^TX}{m} \text{ ולכן עם אותו וקטור עצמי, אזי } L^2 \text{ הוא גם ערך עצמי}$$

$$\text{של } \frac{X^TX}{m} \times \frac{X^TX}{m} \text{ והוא גם מקסימלי עבור מטריצה זו אחרת זו סתירה שהוא מקסימלי עבור המטריצה } \frac{X^TX}{m}.$$

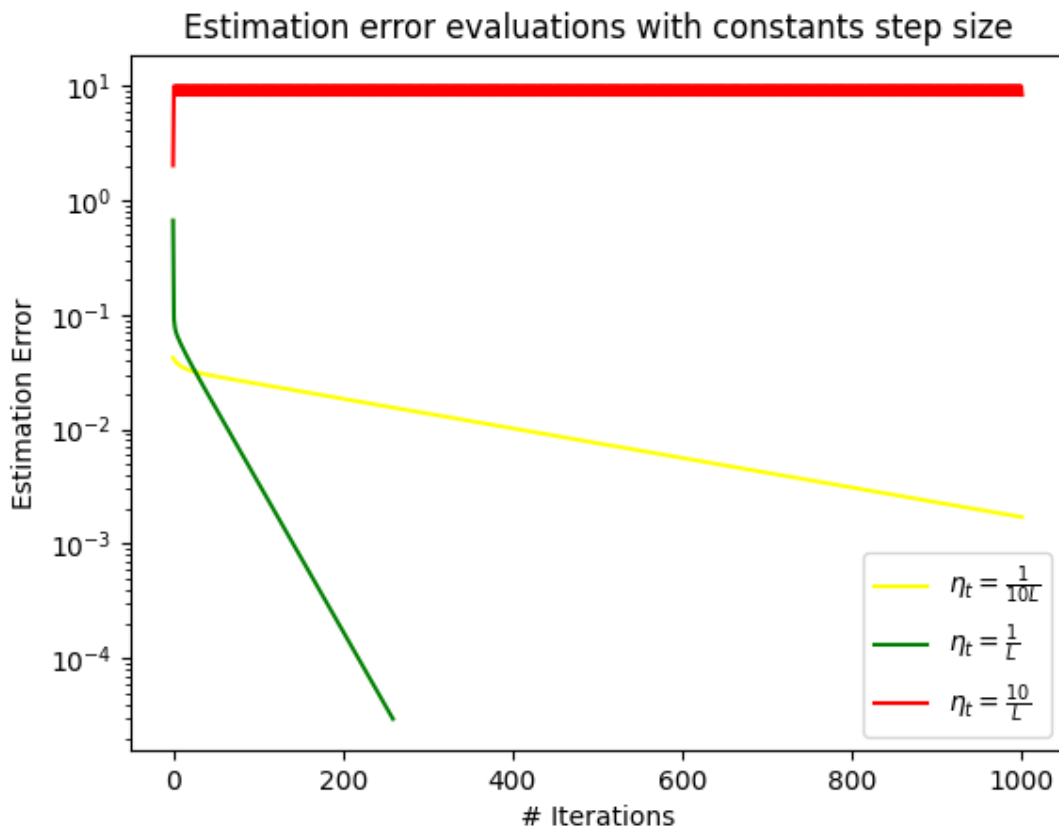
לכן:

$$a_1, a_2 \in C$$

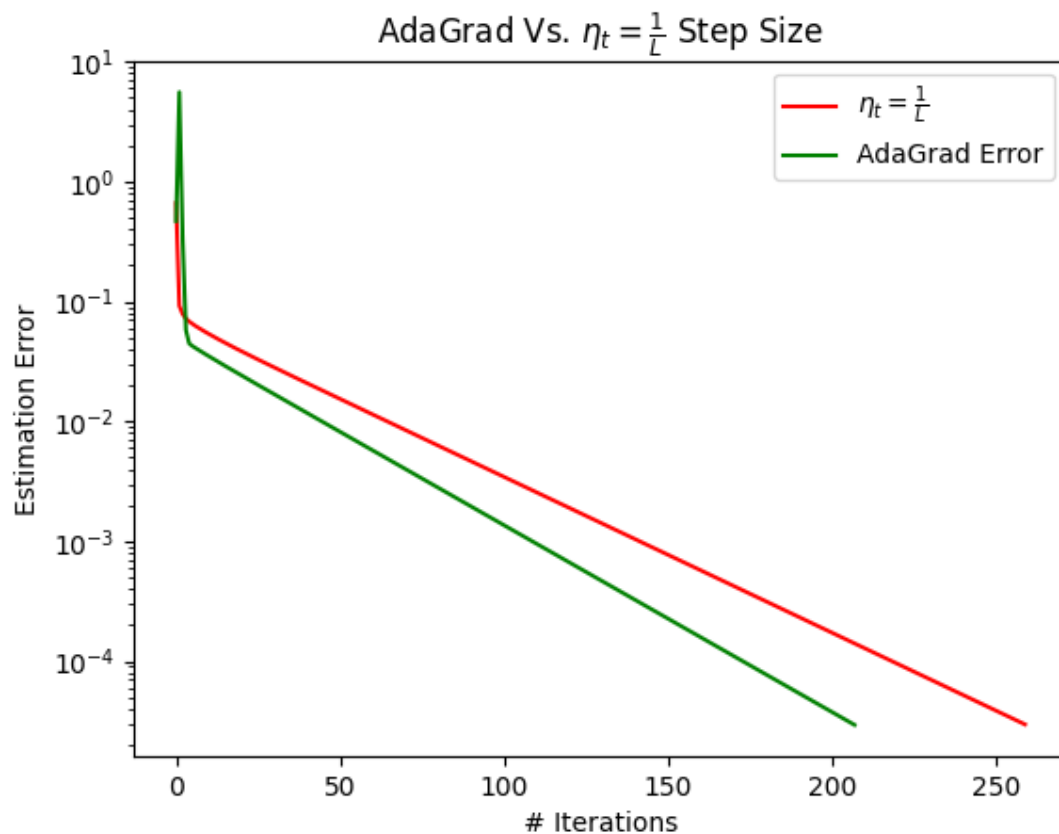
$$\left\|\nabla h(a_1)-\nabla h(a_2)\right\|_2^2 \leq L^2\left\|a_1-a_2\right\|_2^2 \quad .$$

$$\Rightarrow\left\|\nabla h(a_1)-\nabla h(a_2)\right\|_2 \leq L\left\|a_1-a_2\right\|_2$$

ומכיוון שהפונקציה קמורה וגזירה ברציפות ובעל גראדינט לפישצי עם קבוע L אזי היא L חלקה.



מכיוון שאנו מניחים כי דרגת המטריצה X היא $n+1$ הרי הדטרמיננטה שלה שונה מאפס בוודאות ומכיוון ש $\nabla^2 h(a) = \frac{1}{m} (X^T X)$, אזי ההיסאן של פונקציית המטרה בעל דטרמיננטה שונה מאפס בוודאות ולכן בהכרח $\nabla^2 h(a) = \frac{1}{m} (X^T X) \succ 0$ כלומר פונקציית המטרה היא strongly convex ו L חלקה אזי כפי ראינו בכיתה נוכל לבחור בגודל צעד קבוע השווה ל- $\frac{1}{L}$, וקצב ההתכנסות הוא **מעריכי במספר איטרציות**. בנוסף, ניתן לראות בגרף שגודל צעד קטן מידי מוביל להתכנסות איטית יותר וגודל צעד גדול לא מוביל להתכנסות.

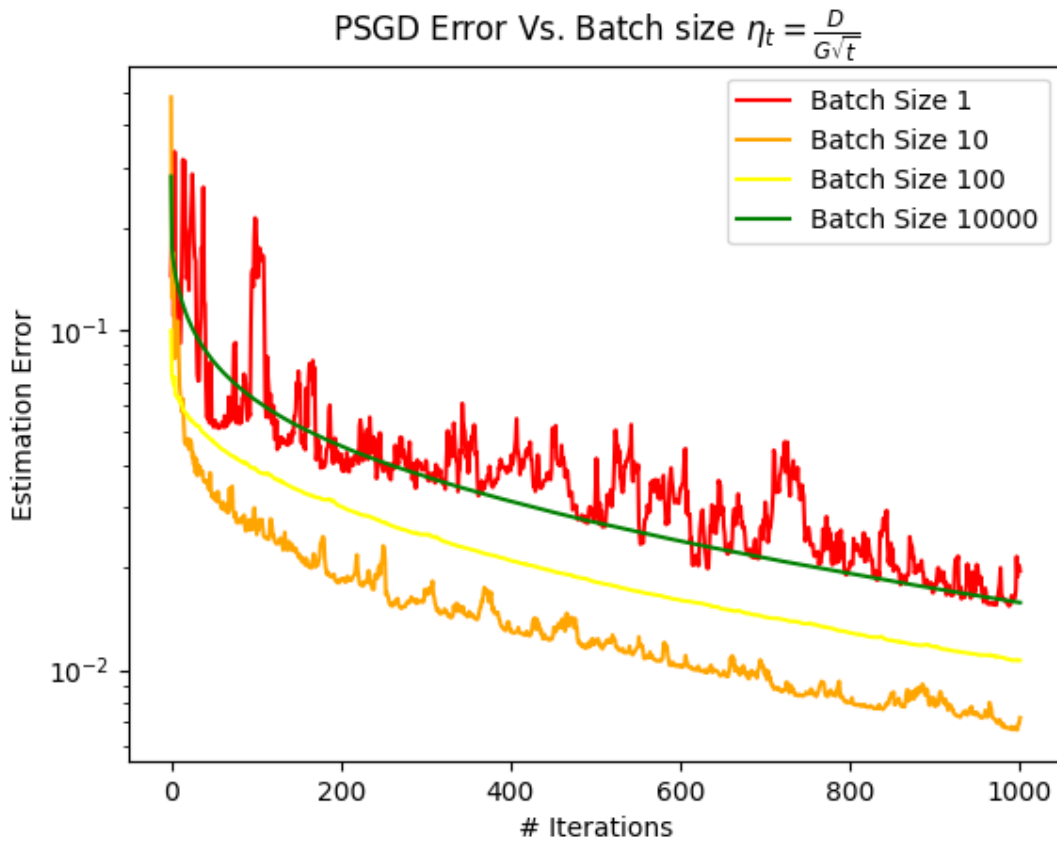


כפי שניתן לראות מהגרף שני גדלי הצעד משיגים קצב התכנסות אקספוננציאלי זאת מכיוון שהגרף הוא לינארי בסקלה לוגריתמית.

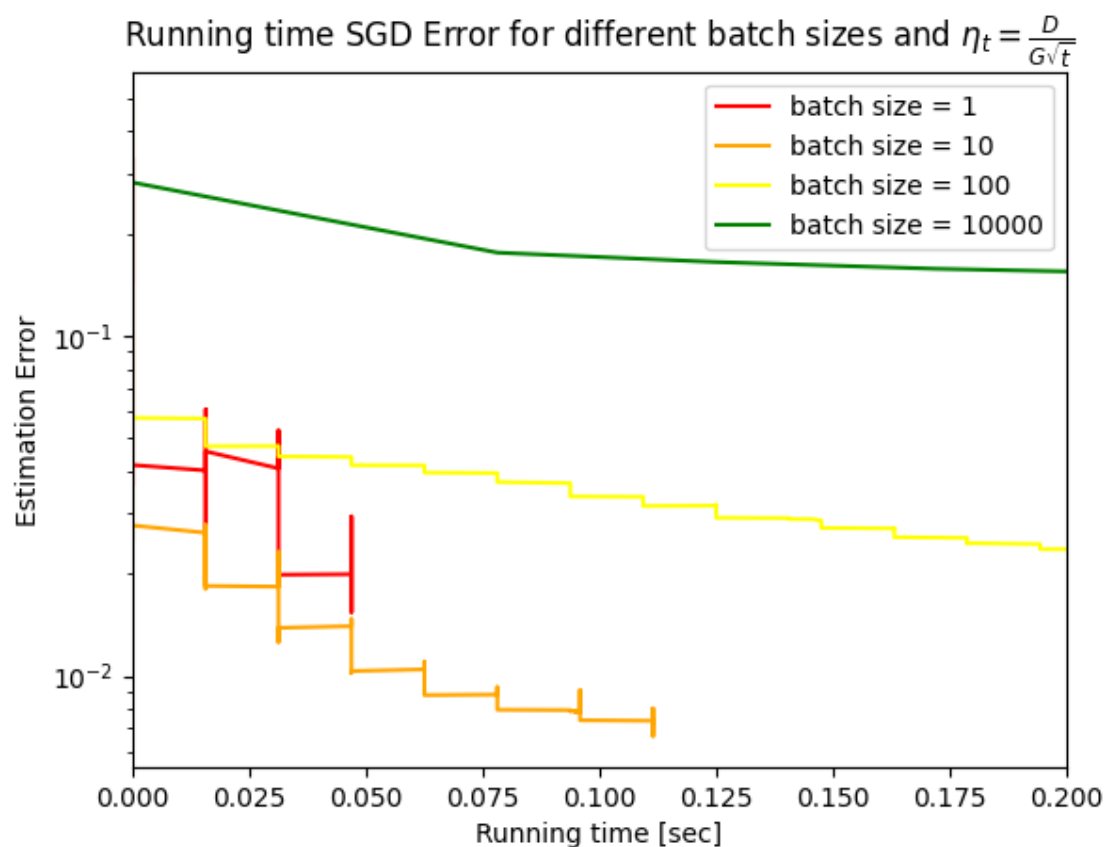
עבור גודל הצעד הקבוע ראינו והוכחנו זאת בתרגול. לעומת זאת, בתרגול הוכחנו כי עבור המקרה

הגרוע ביותר קצב ההתכנסות של גודל צעד *AdaGrad* הוא $O\left(\frac{1}{\sqrt{t}}\right)$. אמנם, מה שמייחד את

AdaGrad שהוא אדפטיבי והוא גורם לגודל הצעד להתאים את עצמו לתהליך עצמו על ידי שימוש בגראדינטים. לכן, הוא עבור מקרה זה גם מושג קצב מעריכי.



ניתן להבין מהגרף כי אלגוריתם ה-SGD הוא יותר רועש מהאלגוריתם ה-GD, זאת מכיוון שהגראדינט נמדד רק על פה כמה מדידות ולא על פי כולן לכן לא בהכרח שבכל איטרציה אנחנו נתקדם לעבר המינימום של הפונקציה הכוללת. כמו שעוד ניתן לשים ככל שאנחנו דוגמים יותר דגימות יחד, כלומר batch גדול יותר התהליך פחות רועש כי כיוון הגראדיאנט יותר נכון. נשים לב, למרות שישנה צפייה כי כאשר אנו דוגמים יותר אז ההתכנסות תהיה מהירה יותר, זה אמנם לא תמיד כך (הדבר משתנה בין הריצות) וזאת מכיוון שההתכנסות תלויה גם בווריאנס של הגראדינט הסטוכסטי.



בגרף הזה אנחנו רואים בבירור את היתרון של SGD ביחס ל-GD. נכון שביחס לכמות איטרציות batch גדול יותר עשוי להתכנס מהר יותר. אמנם, אם נסתכל ביחס ל**זמן** אנחנו רואים כי בכל ש-batch גדול כמו כל הדגימות, אזי השגיאה יורדת בקצב יותר איטי וזאת מכיוון **שכל איטרציה** עבור batch גדול יותר לוקחת יותר זמן (חישוב הגראדינט בכל הדגימות).
 וזה החיסרון של batch גדול יותר, **הזמן בפועל** של ההתכנסות יותר איטי כי כל איטרציה לוקחת יותר זמן, יותר בבדה חישובית, והיתרון הוא שההתכנסות תהיה יותר חלקה ופחות רועשת.