

מערכות לומדות תרגיל 6

עמית בסקין 312259013

PCA

חלק תיאורתי

1. יהא X משתנה מקרי עם התפלגות כלשהי מעל \mathbb{R}^d כך שהתוחלת שלו היא $0_{\mathbb{R}^d}$ ומטריצת השונות המשותפת שלו היא $\Sigma_{d \times d}$.

צ.להוכיח: לכל $v \in \mathbb{R}^d$ עם $\|v\|_2 = 1$, השונות של $\langle v, X \rangle$ אינה גדולה מהשונות של שיכון ה-PCA של X ב- \mathbb{R} .

הוכחה:

לפי ההגדרה מתקיים:

$$\text{Var}(X) = \mathbb{E} \left[(X - \bar{X}) (X - \bar{X})^t \right]$$

כלומר:

$$= \mathbb{E} \left[(X - \bar{0}) (X - \bar{0})^t \right] = \mathbb{E} [X X^t]$$

אז:

$$\text{Var} \langle v, X \rangle = \text{Var} (v^t X) = \mathbb{E} [\langle v^t X, v^t X \rangle] = \mathbb{E} [v^t X X^t v] =$$

$$= v^t \mathbb{E} [X X^t] v = v^t \text{Var}(X) v = v^t \text{Cov}(X) v = v^t \Sigma v$$

יהיו $(x_i)_1^m \subseteq \mathbb{R}^d$ כך ש- $\Sigma = \sum_{i=1}^m x_i x_i^t$ והיו $(u_i)_1^n$ הוקטורים העצמיים המובילים של Σ . נגדיר את U להיות המטריצה שהעמודות

שלה הן $(u_i)_1^n$ ונתבונן בפירוק הספקטרלי של Σ , קרי: $\Sigma = U^t D U$

לפי מה שראינו בהרצאה, מתקיים:

$$\arg \min_{W \in \mathbb{R}^{1 \times d}, U \in \mathbb{R}^{d \times 1}} \sum_{i=1}^m \|x_i - U W x_i\|^2 = (u_1, u_1^t)$$

כלומר, (u_1, u_1^t) הוא פיתרון לבעיית ה-PCA עבור $(x_i)_1^m$.

לפי מה שראינו בהרצאה, מתקיים:

$$\text{trace}(v^t Dv) \leq D_{1,1} = \text{trace}(u_1^t D u_1)$$

אבל:

$$\begin{aligned}\text{trace}(v^t Dv) &= \text{trace}(v^t \Sigma v) = v^t \Sigma v \\ \text{trace}(u_1^t D u_1) &= \text{trace}(u_1^t \Sigma u_1) = u_1^t \Sigma u_1\end{aligned}$$

ולכן:

$$v^t \Sigma v \leq u_1^t \Sigma u_1$$

כאשר $\text{Var}(u_1 X) = u_1^t \Sigma u_1$, קרי השונות של השיכון של X ב- \mathbb{R} . מש"ל

חלק תכנותי

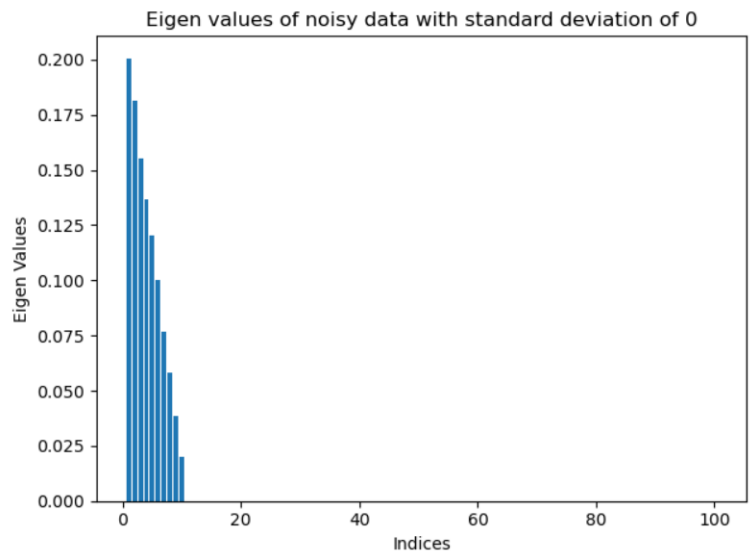
2. נבחן מקרה שבו יש לנו n דגימות מתוך $\mathbb{R}^d \supseteq \mathbb{R}^k$ עבור $k < d$, קרי $d - k$ השיעורים האחרונים של כל דגימה הם אפסים. נניח שהדגימות נדגמו עם רעש גאוסיאני אדיטיבי.

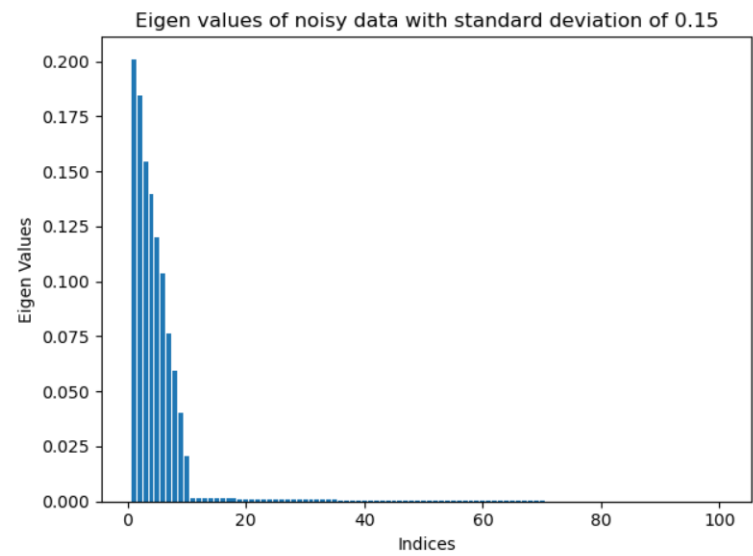
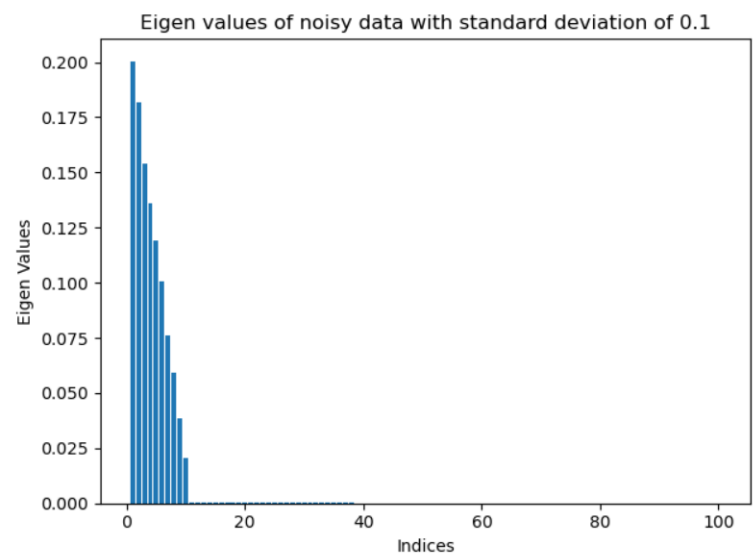
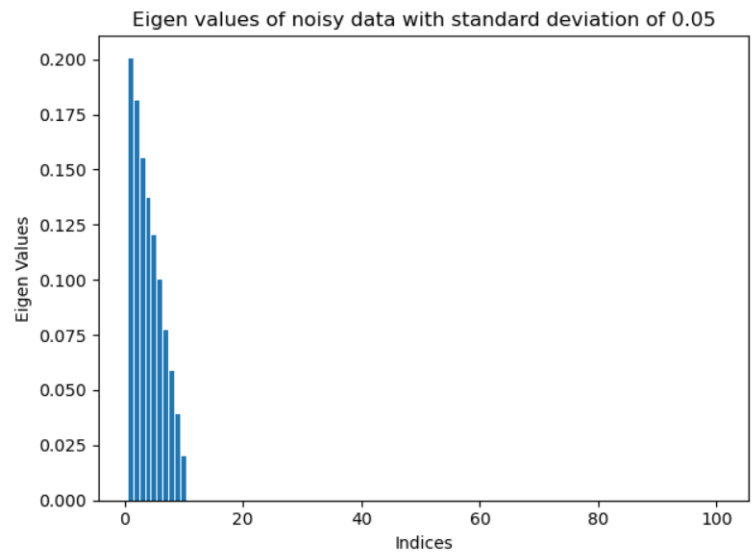
לפני שנטפל בדגימות הרועשות, ננסה להיפטר מהרעש ה"טהור" שקיבלנו ב- $d - k$ הקורדינטות האחרונות של כל דגימה. נראה איך ניתן לבצע זאת באמצעות PCA:

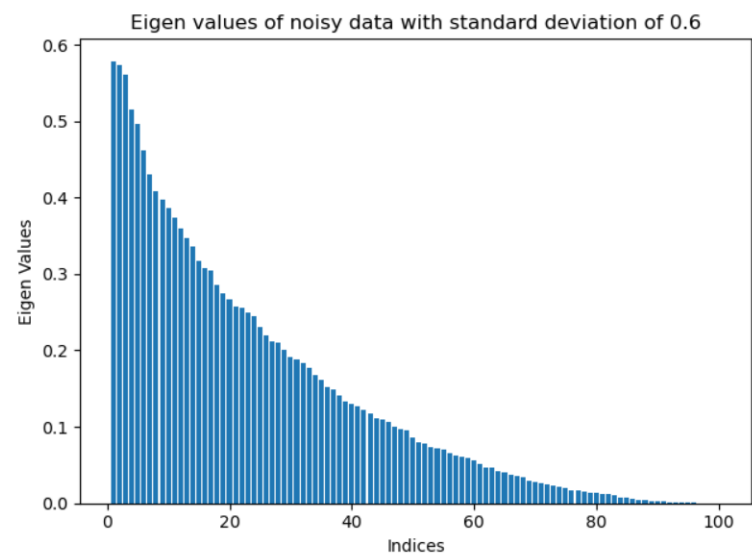
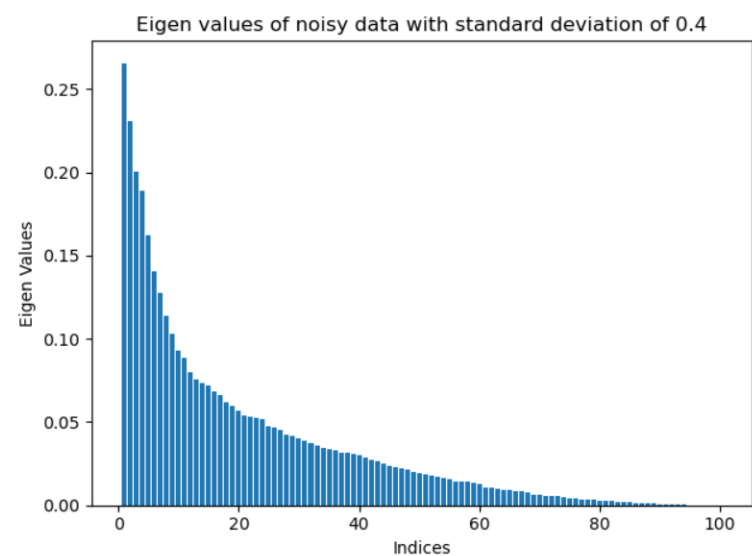
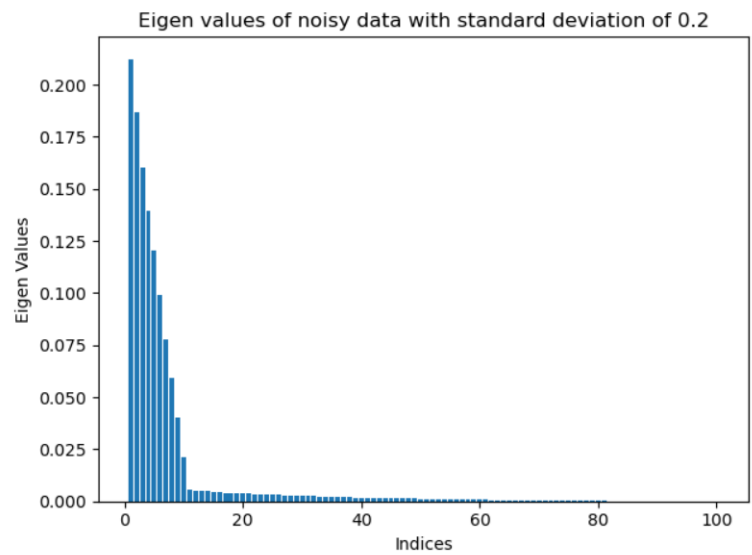
נייצר דאטא $X_{n \times d}$ עם $n = d = 100$ כך ש- X היא מדרגה 10, והערכים הסינגולריים שלה הם $(\sqrt{20 - 2i})_0^9$. בהינתן σ , נדגום מטריצה רנדומלית $Z_{n \times d}$ כאשר כל הכניסות נדגמות באופן זהה ובלתי-לוי מההתפלגות $\mathcal{N}(0, \sigma^2)$. נייצר מטריצת דאטא רועשת: $Y = X + Z$.

נחשב את הערכים העצמיים, $(\lambda_i)_1^{100}$ (בסדר יורד), של מטריצת השונות המשותפת של הדאטא הרועש (כל שורה היא דגימה). נשרטט גרף עמודות של הערכים העצמיים בסדר יורד: ציר ה- x יהיה המספרים הסידוריים של הערכים העצמיים וציר ה- y יהיה הערכים עצמם.

א. נבצע את ה"ל" לכל $\sigma \in \{0, 0.05, 0.1, 0.15, 0.2, 0.4, 0.6\}$:







ב. התופעה שעולה מהנ"ל: ככל שהרעש גדל כך הערכים העצמיים גדלים.

ג. נחקור את הסוגייה של בחירת מימד השיכון:

i. בהינתן הערכים העצמיים של מטריצת השונות המשותפת, צלחשוב על דרך להבדיל בין ערך עצמי שהתקבל כתוצאה

מרעש "טהור" ובין ערך עצמי שהתקבל מהדאטא המקורי.

צ.לתאר את השיטה וליישם אותה עבור הערכים מסעיף א:

כשסטיית התקן קטנה מ- 0.6 נראה שניתן לקבוע באופן די ברור את הרף של הפרש בין שני ערכים עצמיים עוקבים כך

שאם ההפרש קטן מהרף הזה אז מדובר בערכים עצמיים שהתקבלו כתוצאה מרעש טהור.

ניישם את השיטה על הערכים מסעיף א:

```
def determine_original_eigen_vals_amount():
    print()
    for sigma in SIGMAS:
        print('Amount of original eigen values with '
              'standard deviation of {} is '.
              .format(sigma), determine_original_eigen_vals_amount_helper(
                get_eigen_vals(sigma)))

determine_original_eigen_vals_amount()
```

ונקבל:

```
Amount of original eigen values with standard deviation of 0 is 10
Amount of original eigen values with standard deviation of 0.05 is 10
Amount of original eigen values with standard deviation of 0.1 is 12
Amount of original eigen values with standard deviation of 0.15 is 12
Amount of original eigen values with standard deviation of 0.2 is 11
Amount of original eigen values with standard deviation of 0.4 is 11
Amount of original eigen values with standard deviation of 0.6 is 1
```

כלומר, עבור סטיית תקן שקטנה מ- 0.6 השיטה הצליחה פחות או יותר לקבוע מתי מדובר ברעש "טהור" אך עבור 0.6

כבר לא.

ii. צ.לחשב את המכפלה הפנימית שבין עשרת הוקטורים הסינגולריים המובילים של X ובין עשרת הוקטורים העצמיים

המובילים של מטריצת השונות המשותפת של Y :

```
The inner product between the ten leading left singular vectors of X
with the ten leading eigenvectors of the covariance matrix of Y
with standard deviation of 0 is [-0.095-0.014j  0.087-0.005j  0.015+0.11j   0.169-0.033j  0.13 -0.046j
 0.073+0.028j -0.045+0.012j  0.125+0.018j -0.005+0.045j -0.015-0.018j]

The inner product between the ten leading left singular vectors of X
with the ten leading eigenvectors of the covariance matrix of Y
with standard deviation of 0.05 is [ 0.059 -0.045 -0.187 -0.018 -0.118  0.029 -0.277 -0.053  0.075 -0.053]

The inner product between the ten leading left singular vectors of X
with the ten leading eigenvectors of the covariance matrix of Y
with standard deviation of 0.1 is [ 0.122  0.045  0.02  0.059 -0.02  0.002 -0.02 -0.08 -0.007  0.045]

The inner product between the ten leading left singular vectors of X
with the ten leading eigenvectors of the covariance matrix of Y
with standard deviation of 0.15 is [-0.131  0.018 -0.138 -0.031  0.071 -0.079  0.101 -0.092  0.103 -0.185]

The inner product between the ten leading left singular vectors of X
with the ten leading eigenvectors of the covariance matrix of Y
with standard deviation of 0.2 is [-0.008 -0.091  0.137  0.082  0.033 -0.197 -0.04 -0.079 -0.07 -0.078]

The inner product between the ten leading left singular vectors of X
with the ten leading eigenvectors of the covariance matrix of Y
with standard deviation of 0.4 is [-0.087  0.036 -0.201  0.073 -0.058  0.042  0.035 -0.03 -0.12 -0.156]

The inner product between the ten leading left singular vectors of X
with the ten leading eigenvectors of the covariance matrix of Y
with standard deviation of 0.6 is [ 0.04 -0.215 -0.006  0.087 -0.163  0.028 -0.126  0.184 -0.121 -0.166]
```

צ.לתאר את התוצאה שהתקבלה ולהסביר איך זה יכול לעזור לקבוע את השיטה מהסעיף הקודם:

לא ברורה לי התוצאה שהתקבלה אבל אני חושב שהמסקנה שאנו חותרים אליה היא שניתן לחשב את המכפלה הפנימית

של עבור כולם, כלומר לא רק העשרה הראשונים, לפי זה נוכל לומר מתי מדובר ברעש "טהור".

אופטימיזציה קמורה ו- SGD

חלק תיאורתי

3. א. יהיו $\{f_i: V \rightarrow \mathbb{R}\}_1^m$ פונקציות קמורות ו- $\{\gamma_i \in \mathbb{R}_+\}_1^m$. נגדיר: $g: V \rightarrow \mathbb{R}$ על ידי $g(u) = \sum_1^m \gamma_i f_i(u)$. צ'להוכיח ש- g קמורה.

הוכחה:

צ'להראות שלכל $x, y \in V$ ו- $\alpha \in [0, 1]$ מתקיים ש-

$$g(\alpha x + (1 - \alpha)y) \leq \alpha g(x) + (1 - \alpha)g(y)$$

קרי:

$$\begin{aligned} & \sum_1^m \gamma_i f_i(\alpha x + (1 - \alpha)y) \leq \\ & \leq \alpha \sum_1^m \gamma_i f_i(x) + (1 - \alpha) \sum_1^m \gamma_i f_i(y) \end{aligned}$$

ואכן:

$$g(\alpha x + (1 - \alpha)y) =$$

$$= \sum_1^m \gamma_i f_i(\alpha x + (1 - \alpha)y)$$

ומכך שכל f_i קמורה:

$$\begin{aligned} & \leq \sum_1^m \gamma_i (\alpha f_i(x) + (1 - \alpha)f_i(y)) = \\ & = \alpha \sum_1^m \gamma_i f_i(x) + (1 - \alpha) \sum_1^m \gamma_i f_i(y) = \\ & = \alpha g(x) + (1 - \alpha)g(y) \end{aligned}$$

מש"ל.

ב. יהיו $f, g: \mathbb{R} \rightarrow \mathbb{R}$. נגדיר $h: \mathbb{R} \rightarrow \mathbb{R}$ על ידי $h = f \circ g$. צ'להראות שאם f, g קמורות אז h אינה בהכרח קמורה.

הוכחה:

נראה דוגמה ל- $f, g: \mathbb{R} \rightarrow \mathbb{R}$ ו- $x, y \in \mathbb{R}$ ו- $\alpha \in [0, 1]$ כך ש-

$$f(g(\alpha x + (1 - \alpha)y)) > \alpha f(g(x)) + (1 - \alpha)f(g(y))$$

ניקח $g(x) = e^x$ ו- $f(x) = -x$ ו- $x = 2, y = 4$ ו- $\alpha = \frac{1}{2}$. אכן f, g קמורות.

$$f\left(\frac{1}{2} \cdot 2 + \left(1 - \frac{1}{2}\right) \cdot 4\right) =$$

$$= f(g(1+2)) = f(g(3)) = f(e^3) = -e^3 \approx -20$$

ומצד שני:

$$\frac{1}{2}f(g(2)) + \left(1 - \frac{1}{2}\right)f(g(4)) =$$

$$= \frac{1}{2}f(e^2) + \frac{1}{2}f(e^4) =$$

$$= -\frac{1}{2}e^2 - \frac{1}{2}e^4 \approx -31$$

ג. תהא $f: C \rightarrow \mathbb{R}$ עם C קמורה. נגדיר את האפיגרף של f להיות:

$$E = \{(u, t) \mid f(u) \leq t\}$$

צ. להוכיח ש- f קמורה אם האפיגרף של f קמורה.

הוכחה:

כיוון ראשון: נניח ש- f קמורה ונראה ש- E קמורה:

יהיו $(u, t), (v, s) \in E$. תהא $\alpha \in [0, 1]$. צ. להראות ש-

$$\alpha(u, t) + (1 - \alpha)(v, s) =$$

$$= (\alpha u + (1 - \alpha)v, \alpha t + (1 - \alpha)s) \in E$$

קרי ש-

$$f(\alpha u + (1 - \alpha)v) \leq \alpha t + (1 - \alpha)s$$

ואכן מכך ש- f קמורה:

$$f(\alpha u + (1 - \alpha)v) \leq \alpha f(u) + (1 - \alpha)f(v)$$

ומכך ש- $(u, t), (v, s) \in E$:

$$\leq \alpha t + (1 - \alpha)s$$

מש"ל כיוון ראשון.

כיוון שני: נניח ש- E קמורה ונראה ש- f קמורה:

יהיו $x, y \in C$ ו- $\alpha \in [0, 1]$. צ'להוכיח ש-

$$f(\alpha x + (1 - \alpha)y) \leq \alpha f(x) + (1 - \alpha)f(y)$$

תחילה מכך ש- C קמורה, אזי ש- $\alpha x + (1 - \alpha)y \in C$ ולכן $f(\alpha x + (1 - \alpha)y)$ קיימת בכלל.

עתה בפרט מתקיים ש- $f(x) \leq f(x)$ ו- $f(y) \leq f(y)$ ולכן $(x, f(x)), (y, f(y)) \in E$.

אז מכך ש- E קמורה, נקבל שגם:

$$\alpha(x, f(x)) + (1 - \alpha)(y, f(y)) \in E$$

קרי:

$$(\alpha x + (1 - \alpha)y, \alpha f(x) + (1 - \alpha)f(y)) \in E$$

קרי:

$$f(\alpha x + (1 - \alpha)y) \leq \alpha f(x) + (1 - \alpha)f(y)$$

כנדרש.

ד. יהיו $\{f_i: V \rightarrow \mathbb{R}\}_{i \in I}$ קמורות. נגדיר: $f: V \rightarrow \mathbb{R}$ על ידי $f(u) = \sup_{i \in I} f_i(u)$. צ'להוכיח ש- f קמורה.

הוכחה:

$$E = \text{epi}(f)$$

הנחה סמוייה: V קמורה (נובע מההגדרה של פונקציה קמורה, קרי מכך ש- f_i קמורות).

נראה ש- E קמורה ונקבל מהסעיף הקודם ש- f קמורה:

$$E = \{(u, t) \mid f(u) \leq t\} =$$

$$= \left\{ (u, t) \mid \sup_{i \in I} f_i(u) \leq t \right\}$$

יהיו $(u, t), (v, s) \in E$. תהא $\alpha \in [0, 1]$. צ'להראות ש-

$$\alpha(u, t) + (1 - \alpha)(v, s) =$$

$$= (\alpha u + (1 - \alpha)v, \alpha t + (1 - \alpha)s) \in E$$

קרי ש-

$$f(\alpha u + (1 - \alpha)v) \leq \alpha t + (1 - \alpha)s$$

$$\sup_{i \in I} f_i (\alpha u + (1 - \alpha) v) \leq \alpha t + (1 - \alpha) s$$

אבל לכל i מתקיים ש-

$$f_i (\alpha u + (1 - \alpha) v) \leq \alpha t + (1 - \alpha) s$$

$$M = \sup_{i \in I} f_i (\alpha u + (1 - \alpha) v)$$

אם I סופית אז:

$$\sup_{i \in I} f_i (\alpha u + (1 - \alpha) v) = \max_{i \in I} f_i (\alpha u + (1 - \alpha) v)$$

יהא k עס:

$$k \in \arg \max_{i \in I} f_i (\alpha u + (1 - \alpha) v)$$

קרי:

$$\max_{i \in I} f_i (\alpha u + (1 - \alpha) v) =$$

$$= f_k (\alpha u + (1 - \alpha) v)$$

ומכך ש- f_k קמורה:

$$\leq \alpha t + (1 - \alpha) s$$

כפי שרצינו.

אם I לא סופית:

מהגדרת הסופרמום, לכל n טבעי יש f_{i_n} עס:

$$\sup_{i \in I} f_i (\alpha u + (1 - \alpha) v) < f_{i_n} (\alpha u + (1 - \alpha) v) + \frac{1}{n}$$

כלומר:

$$\lim_{n \rightarrow \infty} f_{i_n} (\alpha u + (1 - \alpha) v) = \sup_{i \in I} f_i (\alpha u + (1 - \alpha) v)$$

אז מכך ש- f_{i_n} קמורה לכל i_n , אזי ש-

$$f_{i_n} (\alpha u + (1 - \alpha) v) - \frac{1}{n} < \alpha t + (1 - \alpha) s$$

ולכן:

$$\lim_{n \rightarrow \infty} \left(f_{i_n} (\alpha u + (1 - \alpha) v) - \frac{1}{n} \right) \leq \alpha t + (1 - \alpha) s$$

קרי:

$$\lim_{n \rightarrow \infty} f_{i_n} (\alpha u + (1 - \alpha) v) \leq \alpha t + (1 - \alpha) s$$

קרי:

$$\sup_{i \in I} f_i(\alpha u + (1 - \alpha) v) \leq \alpha t + (1 - \alpha) s$$

כנדרש.

4. יהיו $x \in \mathbb{R}^d$ ו- $y \in \{\pm 1\}$ נגדיר:

$$f(w, b) = l_{x,y}^{hinge}(w, b) = \max\{0, 1 - y \cdot (\langle w, x \rangle + b)\}$$

צ. להראות ש- f קמורה.

הוכחה:

הפונקציה הקבועה 0 היא קמורה וכן הפונקציה $(w, b) \mapsto 1 - y \cdot (\langle w, x \rangle + b)$ היא פונקציה קמורה, כלומר f היא מקסימום של פונקציות קמורות ומהסעיף הקודם מתקבל שהיא קמורה.

ב. צ. למצוא $g \in \partial l_{x,y}^{hinge}(w, b)$

פיתרון:

נפריד למקרים:

אם $\max\{0, 1 - y \cdot (\langle w, x \rangle + b)\} = 0$ אז לפי טענה מהתרגול מתקיים $\partial l_{x,y}^{hinge}(w, b) = \nabla 0(w, b) = \bar{0}$ אבל $\nabla 0(w, b) = \bar{0}$ אז ניקח $g = \bar{0}$

אם $\max\{0, 1 - y \cdot (\langle w, x \rangle + b)\} = 1 - y \cdot (\langle w, x \rangle + b)$ אז לפי טענה מהתרגול מתקיים:

$$h(w, b) = 1 - y \cdot (\langle w, x \rangle + b)$$

$$\nabla h(w, b) = \partial l_{x,y}^{hinge}(w, b)$$

קרי:

$$\begin{pmatrix} -yx_1 \\ -yx_2 \\ \vdots \\ -yx_d \\ -y \end{pmatrix} = -y \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \\ 1 \end{pmatrix} = -y \begin{pmatrix} x \\ 1 \end{pmatrix}$$

$$g = -y \begin{pmatrix} x \\ 1 \end{pmatrix} \text{ כלומר ניקח}$$

מש"ל

ג. יהיו $\{f_i: \mathbb{R}^d \rightarrow \mathbb{R}\}$ קמורות, ו- $\{\xi_i\}_1^m$ כך ש- $\xi_i \in \partial f_i(x)$ לכל $i \in [m]$. נגדיר $f: \mathbb{R}^d \rightarrow \mathbb{R}$ על ידי $f(x) = \sum_1^m f_i(x)$

צ. להראות ש- $\sum_1^m \xi_i \in \partial \sum_1^m f_i(x)$

הוכחה:

מההגדרה של תת-גרדיינט, לכל i מתקיים שלכל $z \in \mathbb{R}^d$ מתקיים:

$$f_i(z) \geq f_i(x) + \langle \xi_i, z - x \rangle$$

ולכן:

$$\sum_1^m f_i(z) \geq \sum_1^m f_i(x) + \left\langle \sum_1^m \xi_i, z - x \right\rangle$$

קרי $\sum_1^m \xi_i \in \partial \sum_1^m f_i(x)$ כנדרש.

ד. יהיו $\{(x_i, y_i)\}_1^m \subseteq \mathbb{R}^d \times \{\pm 1\}$ ו- $\lambda \geq 0$. נגדיר $f: \mathbb{R}^d \times \mathbb{R} \rightarrow \mathbb{R}$ על ידי:

$$f(w, b) = \frac{1}{m} \sum_1^m l_{x,y}^{hinge}(w, b) + \frac{1}{2} \lambda \|w\|^2$$

צ. להראות ש- f קמורה, ולמצוא $\xi \in \partial f(w, b)$.

פיתרון:

מתקיים:

$$f(w, b) = \frac{1}{m} \sum_1^m \max\{0, 1 - y_i \cdot (\langle w, x_i \rangle + b)\} + \frac{1}{2} \lambda \|w\|^2$$

מ-3. ד. מתקיים ש- $\max\{0, 1 - y_i \cdot (\langle w, x_i \rangle + b)\}$ קמורה לכל i .

ר-3. א. מתקיים ש- $\frac{1}{m} \sum_1^m \max\{0, 1 - y_i \cdot (\langle w, x_i \rangle + b)\}$ קמורה ומכך ש- $\frac{1}{2} \lambda \|w\|^2$ קמורה, אזי שגם $\frac{1}{m} \sum_1^m \max\{0, 1 - y_i \cdot (\langle w, x_i \rangle + b)\} + \frac{1}{2} \lambda \|w\|^2$ קמורה.

נגדיר:

$$h(w, b) = \frac{1}{2} \lambda \|w\|^2$$

אז:

$$\nabla h(w, b) = \begin{pmatrix} \lambda w \\ 0 \end{pmatrix}$$

לכל i נגדיר:

$$g_i(w, b) = \begin{cases} \bar{0} & \text{if } \max\{0, 1 - y_i \cdot (\langle w, x_i \rangle + b)\} = 0 \\ -y_i \begin{pmatrix} x_i \\ 1 \end{pmatrix} & \text{else} \end{cases}$$

אז משני הסעיפים הקודמים מתקבל ש-

$$\begin{pmatrix} \lambda w \\ 0 \end{pmatrix} + \frac{1}{m} \sum_1^m g_i(w, b) \in \partial f(w, b)$$

קרי:

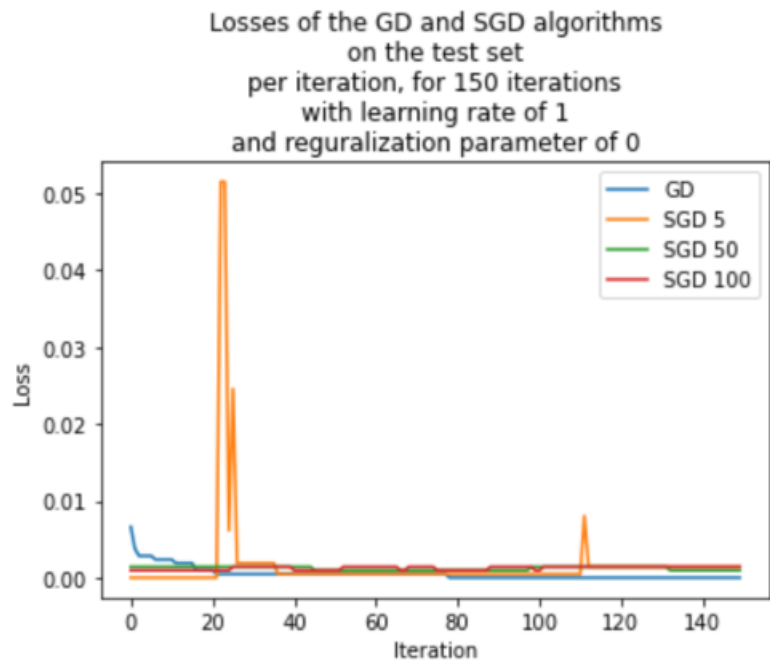
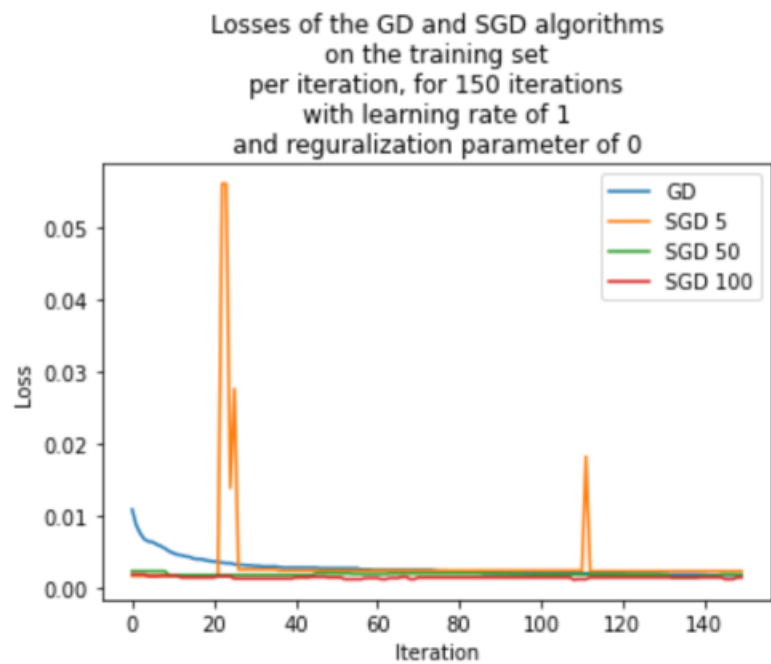
$$\begin{pmatrix} \lambda w \\ 0 \end{pmatrix} - \frac{1}{m} \sum_{y_i (\langle w, x_i \rangle + b) < 1} y_i \begin{pmatrix} x_i \\ 1 \end{pmatrix} \in \partial f(w, b)$$

מש"ל

חלק תכנותי - SVM על MNIST

בחלק זה נפתור בעיית קלאסיפיקציה על הדאטא של MNIST - תמונות של 28×28 פיקסלים של ספרות שנכתבו בכתב יד. א. השתמשתי ב- $\eta = 1$.

ב. ש.לשרטט גרף של ה- $loss$ של האימון עבור כל אחד מהאלגוריתמים בכל איטרציה, וכנ"ל עבור ה- $loss$ של המבחן:



ג. SGD משתמש במיני באטצ'ים ומנסה לשערך את GD אבל עם פחות דגימות.

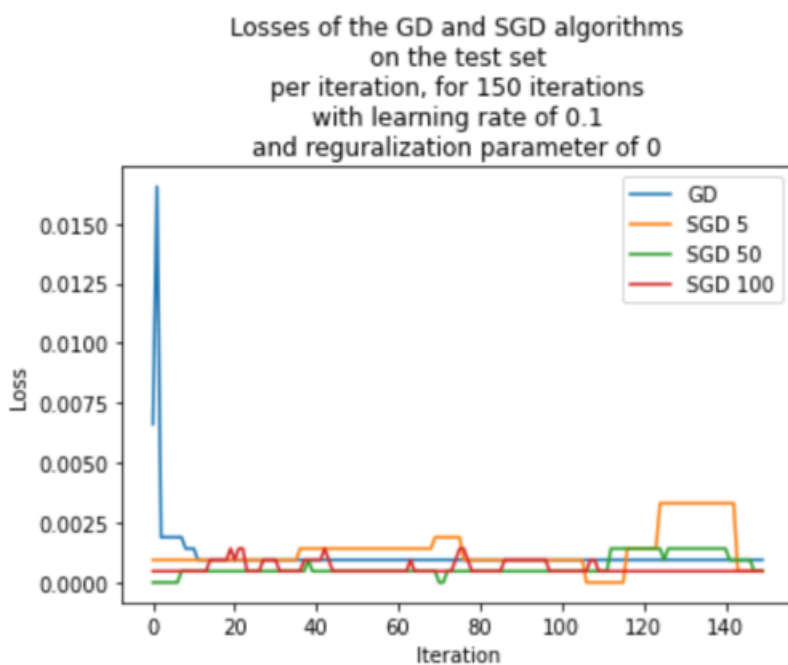
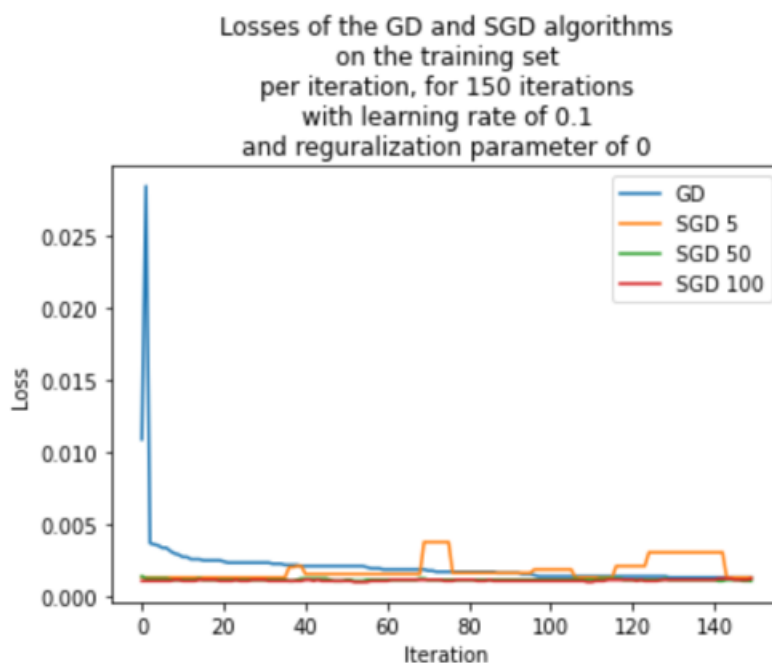
צ. לנתח את התוצאות מהסעיף הקודם. צ. להסביר כיצד השינוי בגודל הבאטצ'ים משפיע על הסיבוכיות החישובית ועל הדיוק של תת-הגרדיינט בתהליך הלמידה.

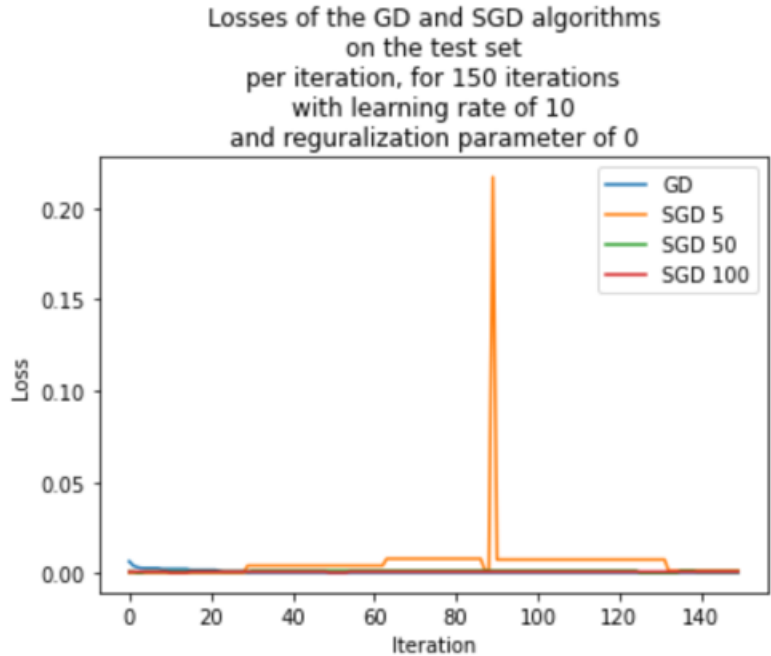
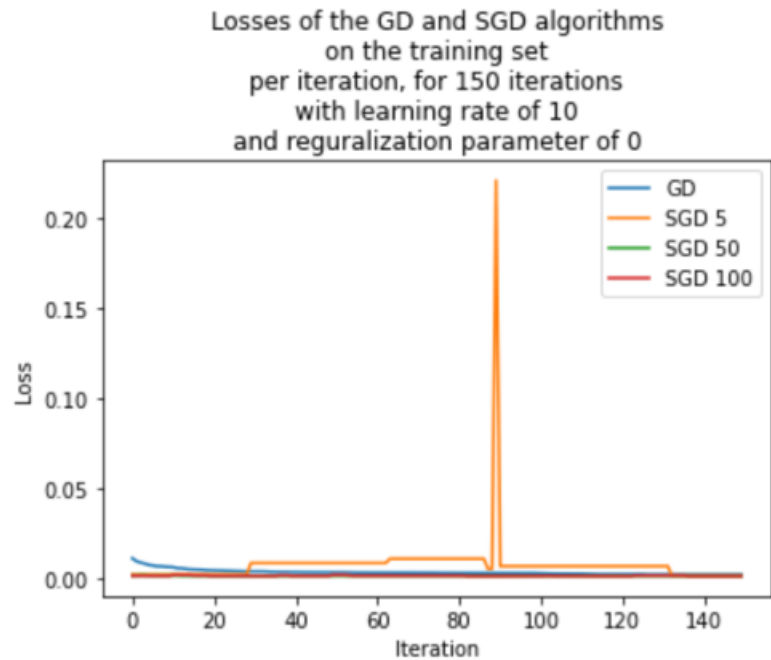
תשובה:

גודל הבאטצ' משפיע בעיקר בחישוב תת-הגרדיינט בכל איטרציה. כל חישוב כזה הוא לינארי בגודל הבאטצ'. נסמן ב- T את מספר האיטרציות וב- B גודל של באטצ' נתון. אז הסיבוכיות החישובית היא TB . לכן, עבור סקלר α כלשהו, אם לוקחים באטצ' מגודל αB , ההבדל בסיבוכיות החישובית הוא $|1 - \alpha|TB$.

באשר לדיוק, ככל השגודל הבאטצ' גדול יותר, כך הוא יותר מייצג את הדאטא הנתון ולכן הדיוק גדל.

ד. צ. לחזור על סעיפים א ו- ב עם $\eta = 0.1, 10$ ולבדוק מה ניתן ללמוד משיעור הלמידה במקרים הללו:





אז במקרה הזה, עבור מקדם למידה גדול יחסית, ה- $loss$ קטן בצורה אחידה יותר על פני האיטרציות.

שאלה: עבור אילו מקרים כדאי לנו לקחת מקדם למידה גבוה / שיעור למידה נמוך? בשביל לענות על השאלה ניתן להשתמש בגרדיינט

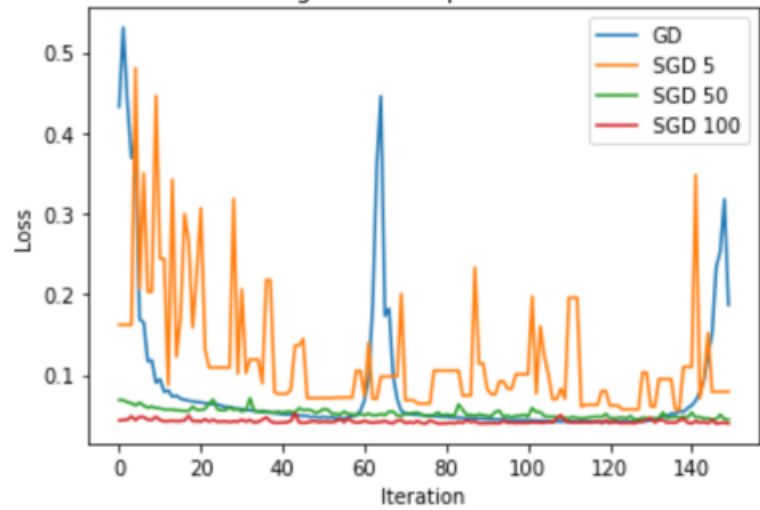
של $f(x) = 10^8 \cdot |x|$ ושל $g(x) = 10^{-8} \cdot |x|$ בנקודות שונות, ובהשפעה שלו ב- GD .

תשובה:

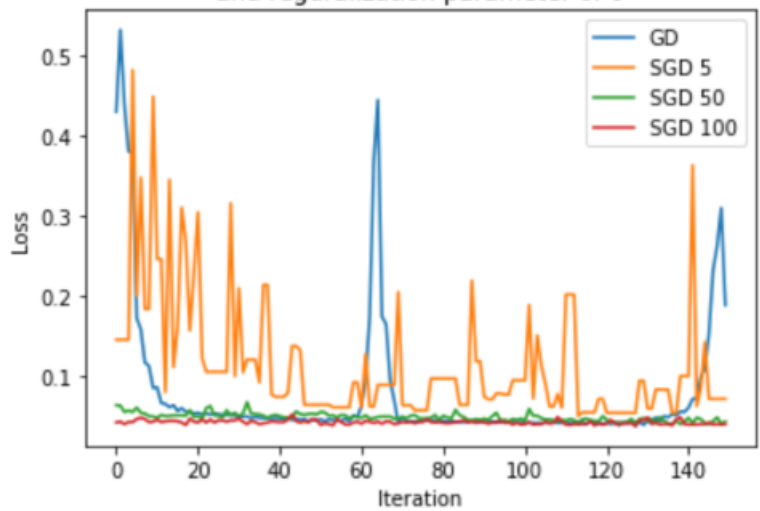
ככל שהת-גרדיינט קטן יותר בערך מוחלט, כך כדאי לקחת מקדם למידה גבוה יותר, על מנת ש"נספיק" למצוא את המינימום הגלובלי, שהרי הפונקציה שעבורה לקחנו תת-גרדיינט, מתקדמת לאט. ובאופן אנלוגי, ככל שהת-גרדיינט גבוה יותר בערך מוחלט, כך כדאי לקחת מקדם למידה נמוך יותר, על מנת שלא "נפספס" את המינימום הגלובלי, שהרי הפונקציה שעבורה לקחנו תת-גרדיינט, מתקדמת מהר.

ה. צלבצע קלאסיפיקציה עבור זוג ספרות נוסף: למצוא זוג ספרות כך שה- GD או ה- SGD לא מתפקד היטב עבורו, ולהציע הסבר עבור התוצאות הגרועות:

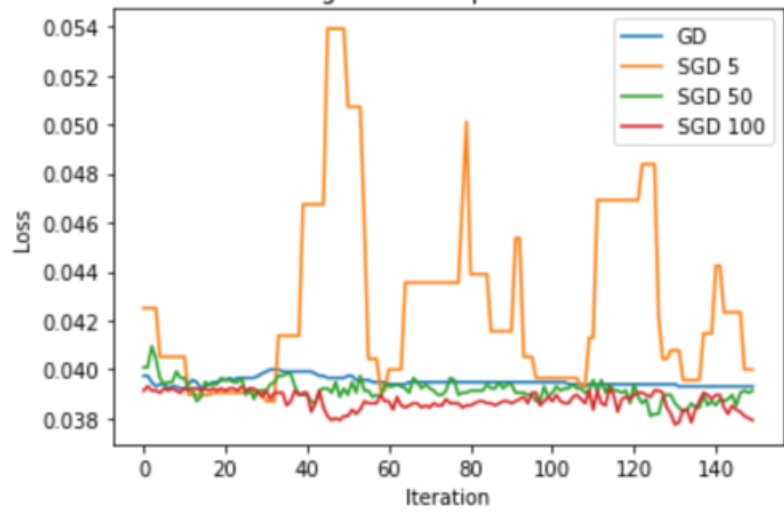
Losses of the GD and SGD algorithms
on the training set
per iteration, for 150 iterations
with learning rate of 1
and reguralization parameter of 0



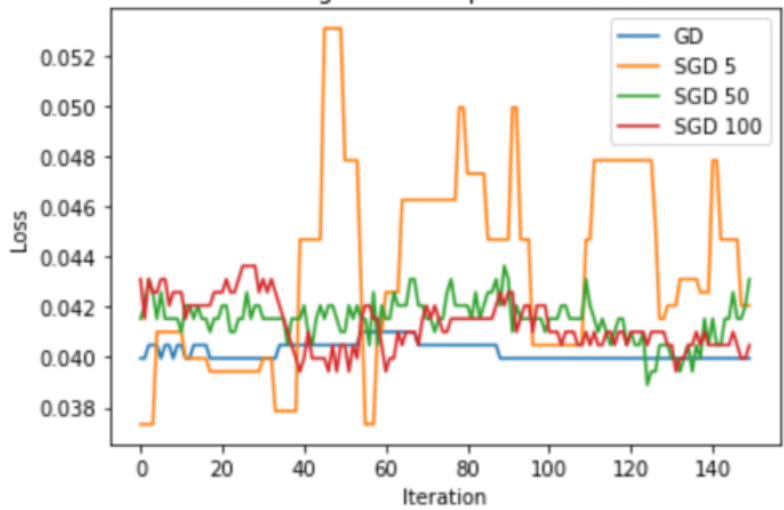
Losses of the GD and SGD algorithms
on the test set
per iteration, for 150 iterations
with learning rate of 1
and reguralization parameter of 0



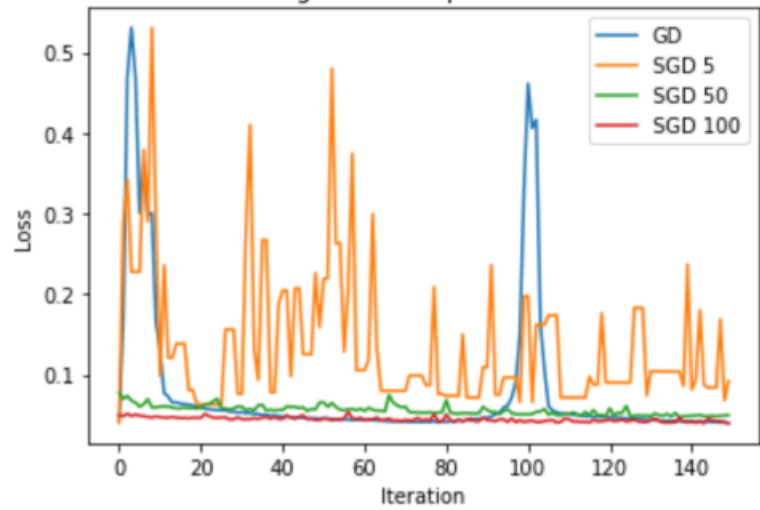
Losses of the GD and SGD algorithms
on the training set
per iteration, for 150 iterations
with learning rate of 0.1
and regularization parameter of 0



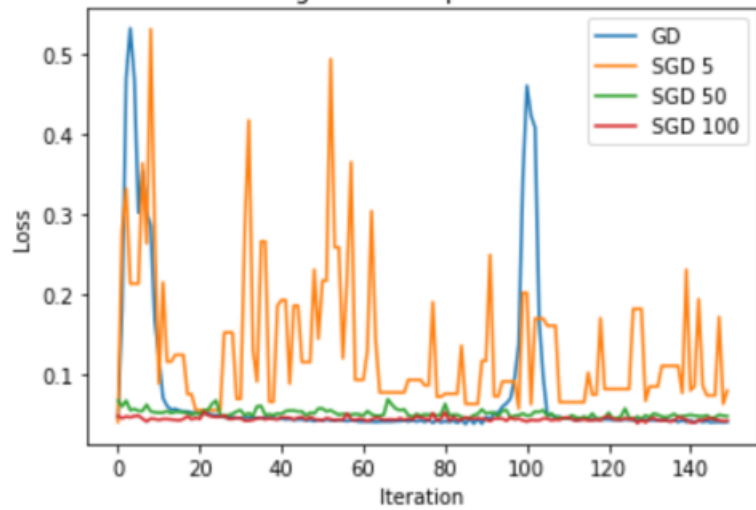
Losses of the GD and SGD algorithms
on the test set
per iteration, for 150 iterations
with learning rate of 0.1
and regularization parameter of 0



Losses of the GD and SGD algorithms
on the training set
per iteration, for 150 iterations
with learning rate of 10
and regularization parameter of 0



Losses of the GD and SGD algorithms
on the test set
per iteration, for 150 iterations
with learning rate of 10
and regularization parameter of 0



הסבר:

הדמיון בכתיב של 3 ו- 5 גדול יבאופן משמעותי מאשר הדמיון שבין הכתיב של 0 ו- 1 ולכן יותר קשה להם להבחין בין הספרות.