

Contents

1	Ex6 Answers.pdf	2
2	PCA.py	19
3	SVM.py	21
4	ex6 runme.py	25

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:

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

כלומר:

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

אז:

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

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

יהיו $(x_i)_1^m \subseteq \mathbb{R}^d$ כך ש- $\Sigma = \sum_1^m x_i x_i^t$ 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:

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

אבל:

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

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 $\mathbb{R}^k \subseteq \mathbb{R}^d$ 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:

נניח דאטא $X_{n \times d}$ עם $n = d = 100$ כך ש- X היא מדרגה 10, והערכים הסינגולריים שלה הם $(\sqrt{20 - 2i})_0^9$.

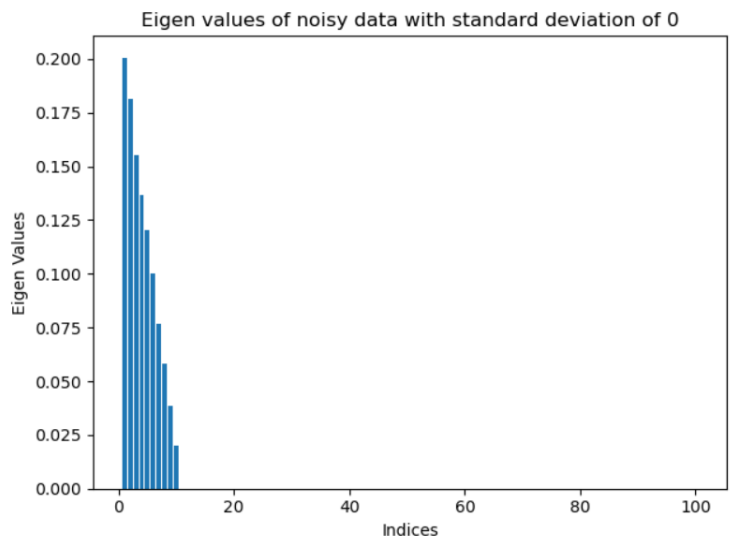
בהינתן σ , נדגום מטריצה רנדומלית $Z_{n \times d}$ כאשר כל הכניסות נדגמות באופן זהה ובלתי-תלוי מההתפלגות $\mathcal{N}(0, \sigma^2)$.

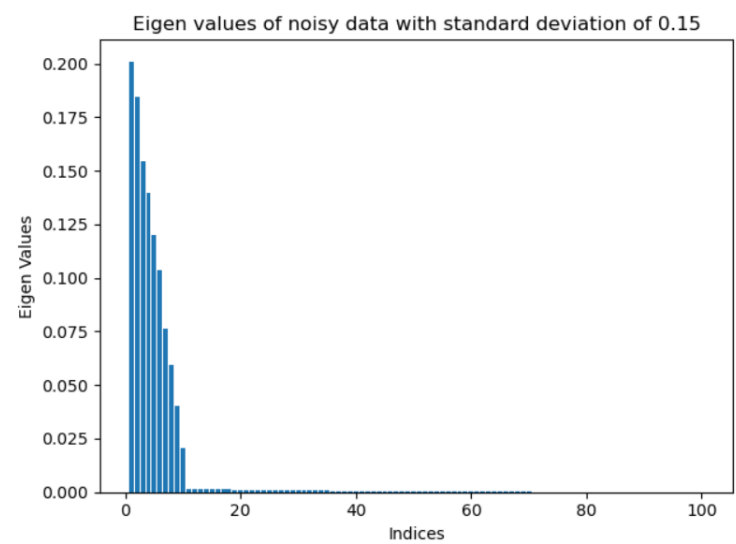
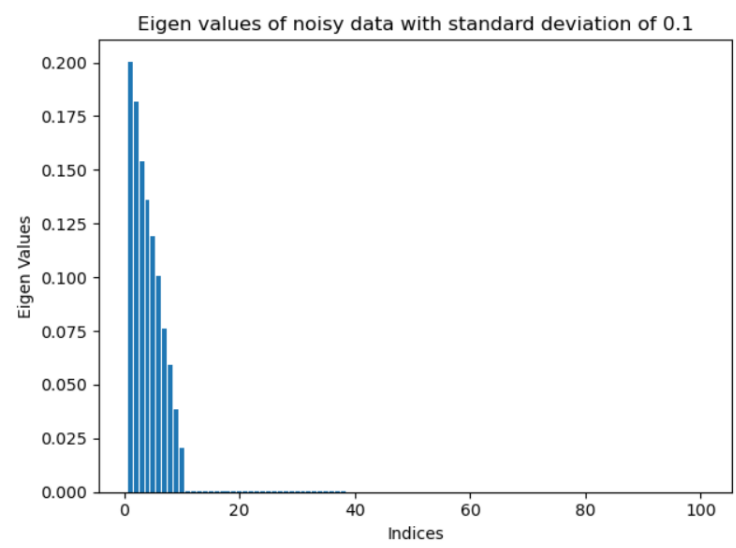
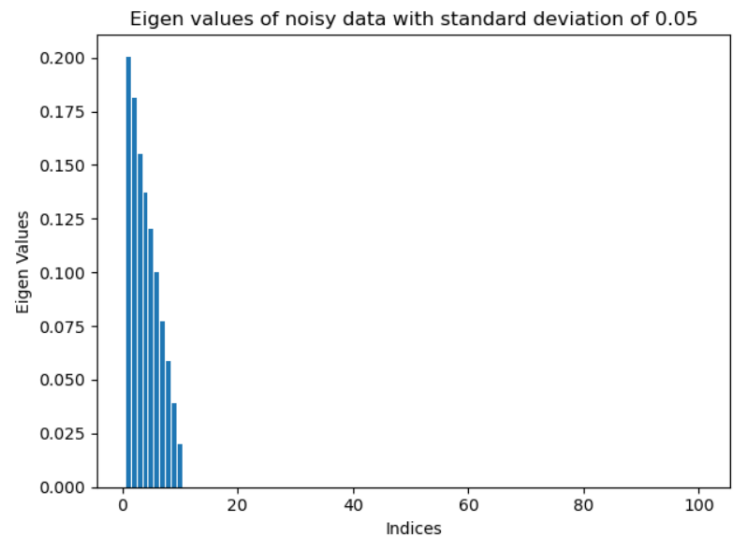
נניח מטריצת דאטא רועשת: $Y = X + Z$.

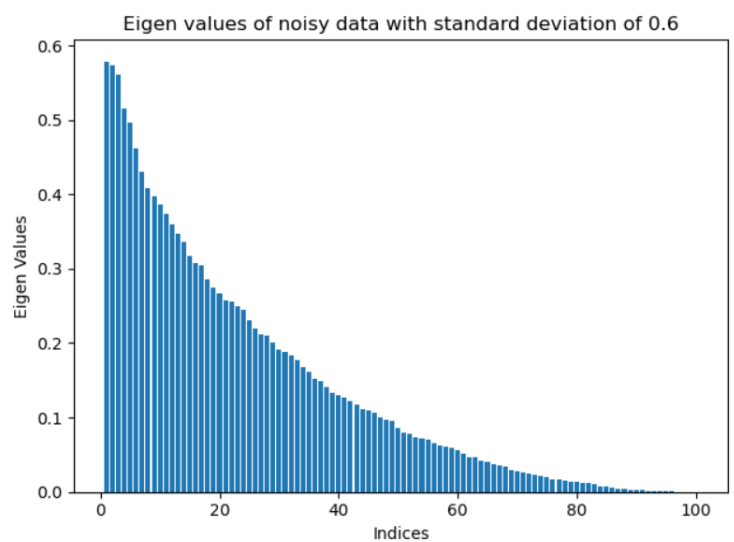
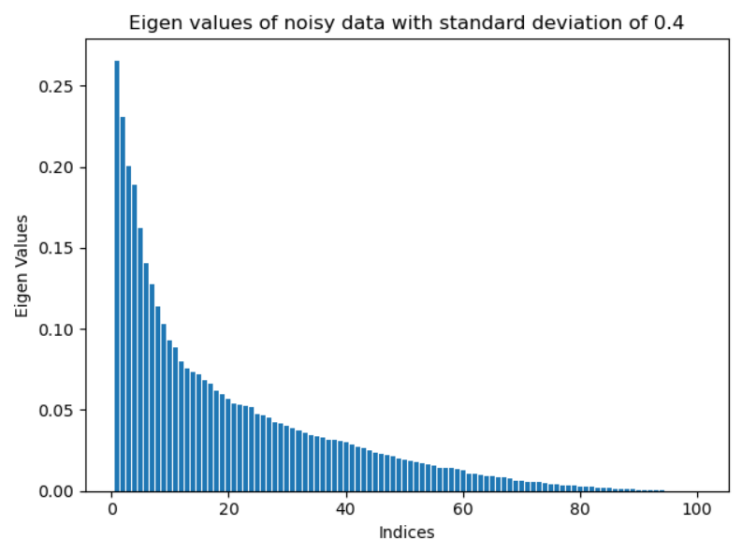
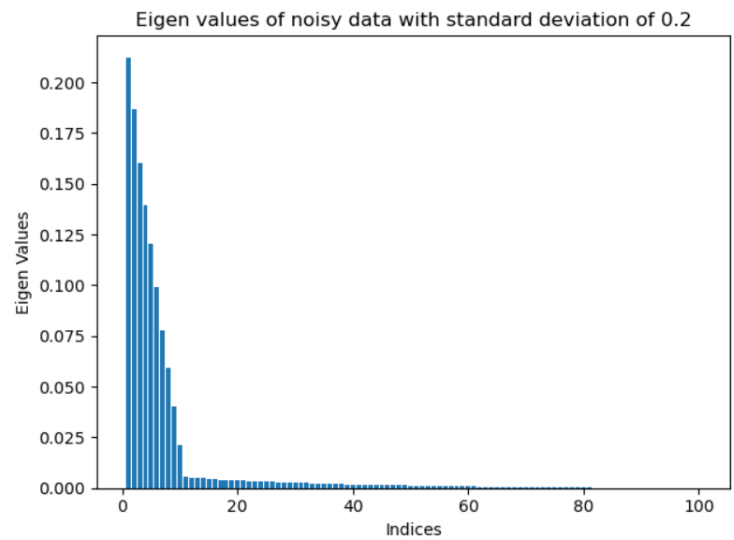
נחשב את הערכים העצמיים, $(\lambda_i)_1^{100}$ (בסדר יורד), של מטריצת השונות המשותפת של הדאטא הרועש (כל שורה היא דגימה).

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.

אז מצד אחד:

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

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

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- $\alpha \in [0, 1]$. צלהוכיח 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- $f(x) \leq f(y)$ 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 $f(u) = \sup_{i \in I} f_i(u)$. צלהוכיח ש- f קמורה.

הוכחה:

נסמן $E = \text{epi}(f)$.

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 $\alpha \in [0, 1]$. צלהראות 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 עם:

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

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} עם:

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

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:

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

צלהראות 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. צלמצוא $g \in \partial l_{x,y}^{hinge}(w, b)$.

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:

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

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

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 $f(x) = \sum_1^m f_i(x)$.

צלהראות 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:

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

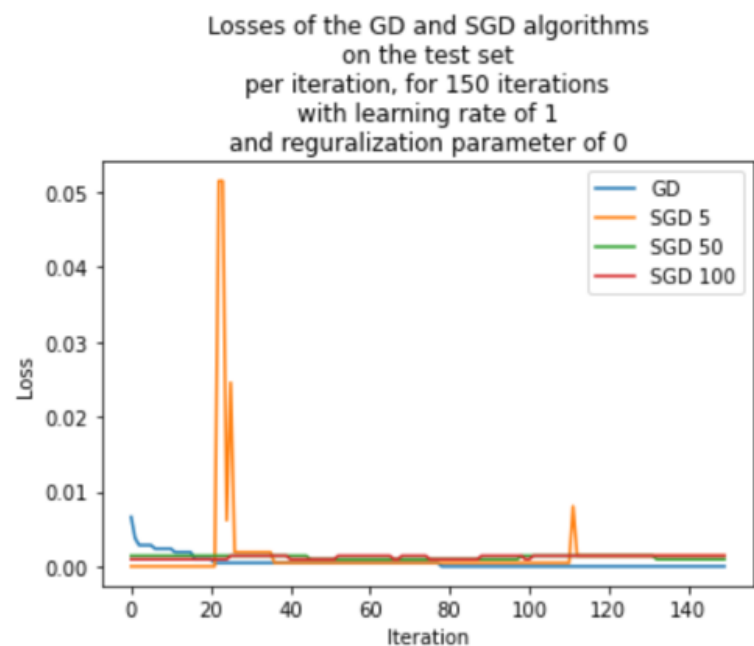
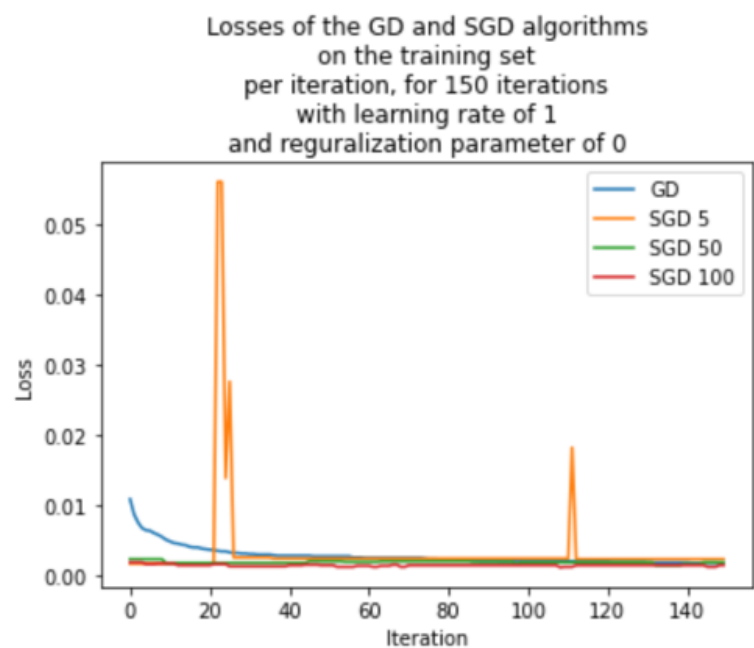
מש"ל

חלק תכנותי - *MNIST* על *SVM*

בחלק זה נפתור בעיית קלאסיפיקציה על הדאטא של *MNIST* - תמונות של 28×28 פיקסלים של ספרות שנכתבו בכתב יד.

6. 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.

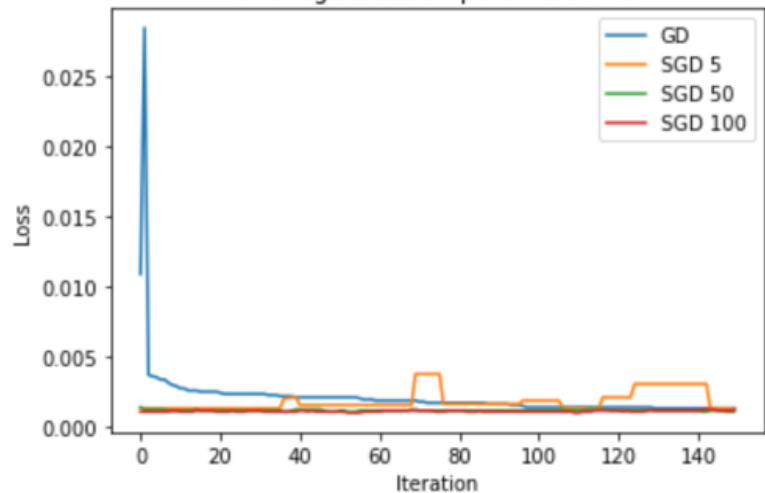
תשובה:

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 $TB|1 - \alpha|$.

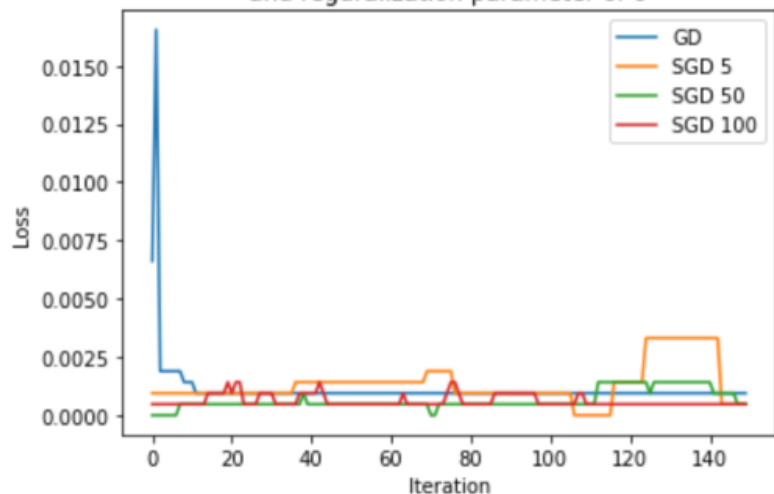
באשר לדיוק, ככל שהגודל הבאטצ' גדול יותר, כך הוא יותר מייצג את הדאטא הנתון ולכן הדיוק גדל.

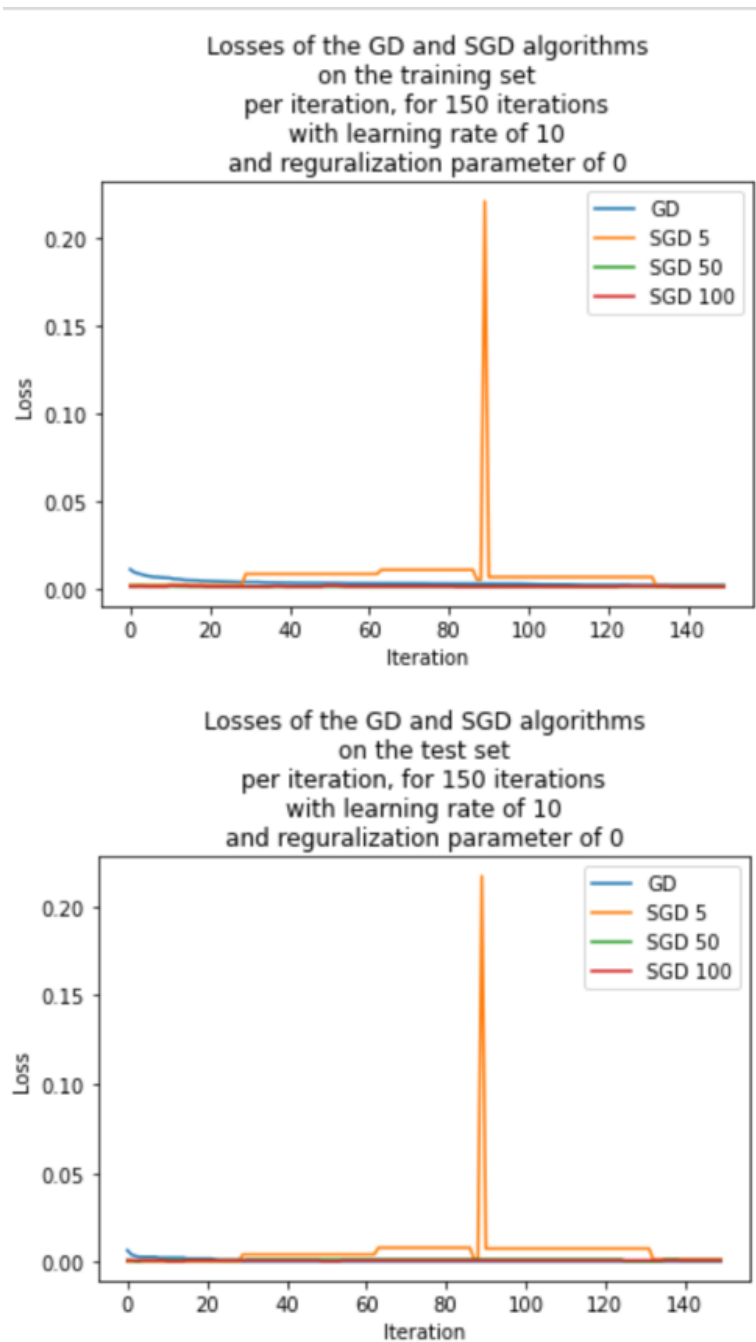
ד. צ. לחזור על סעיפים א ו- ב עם $\eta = 0.1, 10$ ולבדוק מה ניתן ללמוד משיעור הלמידה במקרים הללו:

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





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

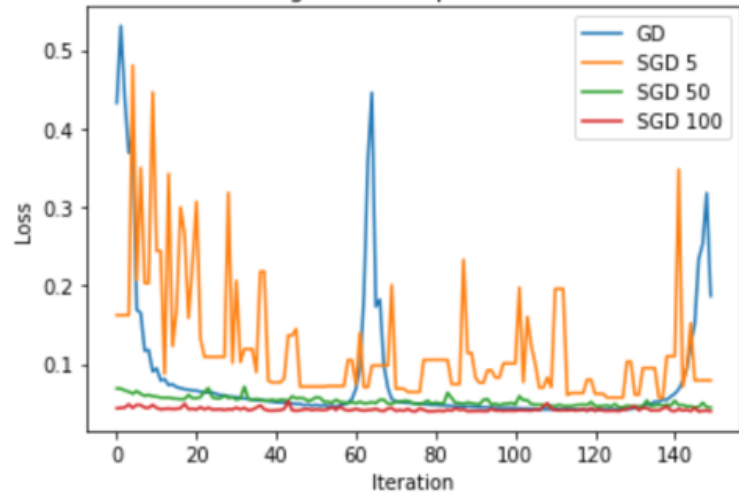
של $f(x) = 10^8 \cdot |x|$ ושל $g(x) = 10^{-8} \cdot |x|$ בנקודות שונות, ובהשפעה שלו ב- GD .

תשובה:

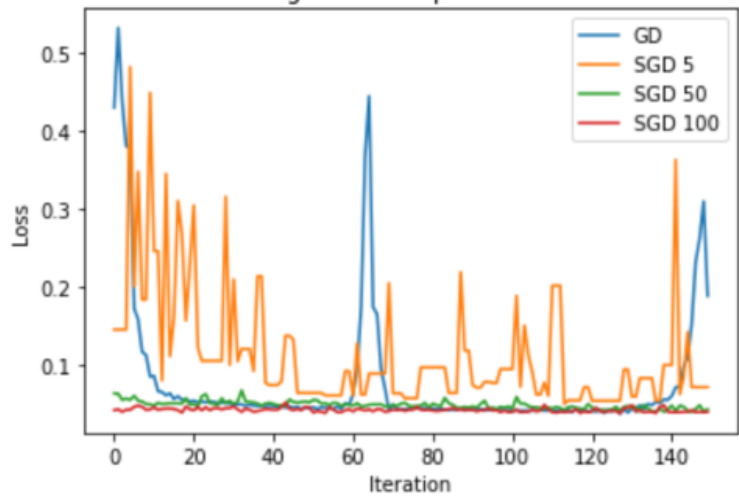
ככל שהת-גרדיינט קטן יותר בערך מוחלט, כך כדאי לקחת מקדם למידה גבוהה יותר, על מנת ש"נספיק" למצוא את המינימום הגלובלי, שהרי הפונקציה שעבורה לקחנו תת-גרדיינט, מתקדמת לאט. ובאופן אנלוגי, ככל שהת-גרדיינט גבוה יותר בערך מוחלט, כך כדאי לקחת מקדם למידה נמוך יותר, על מנת שלא "נפספס" את המינימום הגלובלי, שהרי הפונקציה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

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

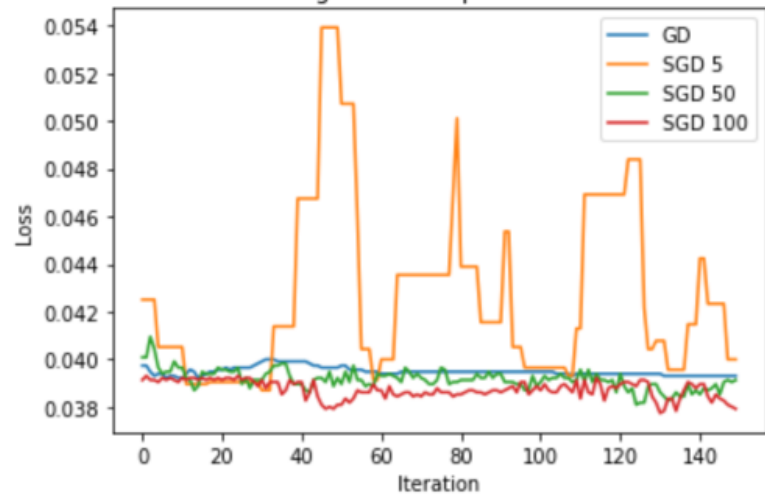
Losses of the GD and SGD algorithms
on the training set
per iteration, for 150 iterations
with learning rate of 1
and regularization parameter of 0



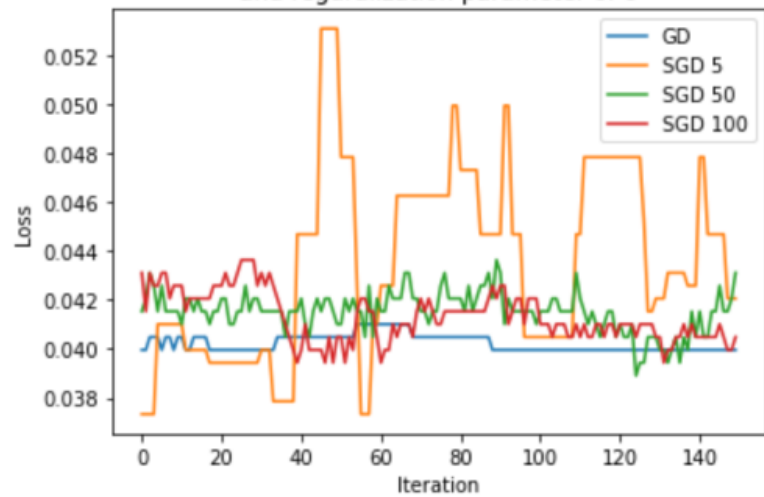
Losses of the GD and SGD algorithms
on the test set
per iteration, for 150 iterations
with learning rate of 1
and regularization parameter of 0



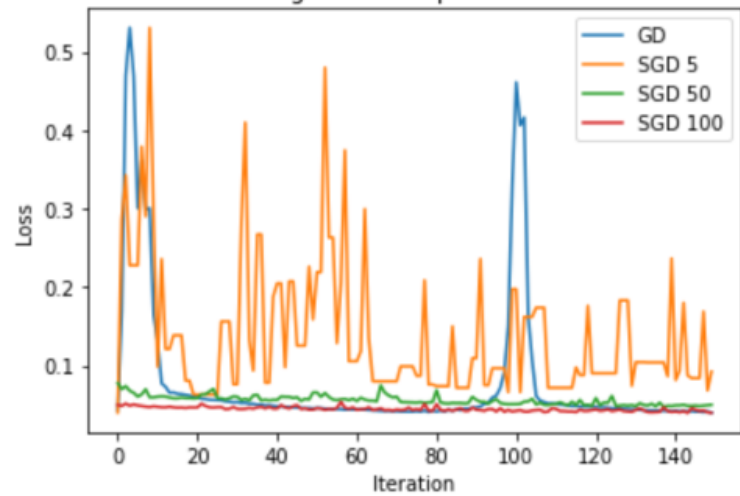
Losses of the GD and SGD algorithms
on the training set
per iteration, for 150 iterations
with learning rate of 0.1
and regularization parameter of 0



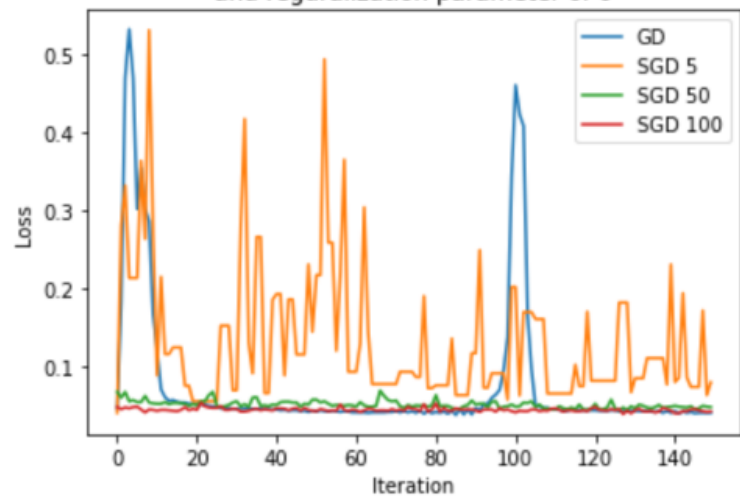
Losses of the GD and SGD algorithms
on the test set
per iteration, for 150 iterations
with learning rate of 0.1
and regularization parameter of 0



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

2 PCA.py

```
1  import numpy as np
2  import matplotlib.pyplot as plt
3
4
5  ROWS_DIM = 100
6  COLS_DIM = 100
7  SIGMAS = [0, 0.05, 0.1, 0.15, 0.2, 0.4, 0.6]
8
9
10 def get_samples():
11     A = np.random.randn(ROWS_DIM, COLS_DIM)
12     A_mean_mat = np.tile(np.mean(A, 0).reshape((1, -1)), (ROWS_DIM, 1))
13     A -= A_mean_mat
14     U, _, V = np.linalg.svd(A)
15     D = np.sqrt(2 * np.arange(10, 0, -1))
16     X = np.dot(np.dot(U[:, :10], np.diag(D)), V[:, 10, :])
17     return X
18
19
20 X = get_samples()
21
22
23 def get_noised_samples(sigma):
24     Z = np.random.normal(0, sigma ** 2, (ROWS_DIM, COLS_DIM))
25     return X + Z
26
27
28 def get_eigenvals(sigma):
29     Y = get_noised_samples(sigma)
30     return np.flip(np.linalg.eigvalsh(np.cov(Y)))
31
32
33 def plot_eigenvals(sigma):
34     eigenvals = get_eigenvals(sigma)
35     indices = np.arange(len(eigenvals)) + 1
36     plt.bar(indices, eigenvals)
37     plt.xlabel('Indices')
38     plt.ylabel('Eigenvalues')
39     plt.title('Eigenvalues of noisy data with standard deviation of {}'.
40             format(sigma))
41     plt.show()
42
43
44 def plot_all_eigenvals():
45     for sigma in SIGMAS:
46         plot_eigenvals(sigma)
47
48
49 # plot_all_eigenvals()
50
51
52 def determine_original_eigenvals_amount_helper(eigenvals):
53     counter = 1
54     for i in range(len(eigenvals) - 1):
55         curr = eigenvals[i]
56         next = eigenvals[i + 1]
57         if 0.95 * curr > next and next > 0.01 * curr:
58             counter += 1
59     else:
```

```

60         break
61     return counter
62
63
64 def determine_original_eigenvals_amount():
65     print()
66     for sigma in SIGMAS:
67         print('    Amount of original eigenvalues with '
68               'standard deviation of {} is '
69               .format(sigma), determine_original_eigenvals_amount_helper(
70                 get_eigenvals(sigma)))
71
72
73 # determine_original_eigenvals_amount()
74
75
76 def get_original_and_noised_product_helper(sigma):
77     Y = get_noised_samples(sigma)
78     U, _, V = np.linalg.svd(X)
79     X_left_singular_vecs = U[:10]
80     Y_cov_eigenvecs = np.linalg.eig(np.cov(Y))[1][:10]
81     inner_product = []
82     for i in range(10):
83         sing = X_left_singular_vecs[i]
84         eig = Y_cov_eigenvecs[i]
85         inner_product.append(sing @ eig)
86     return inner_product
87
88
89 def get_original_and_noised_product():
90     print()
91     for sigma in SIGMAS:
92         print()
93         print('The inner product between the ten leading left singular '
94               'vectors of X\nwith the ten leading eigenvectors of the '
95               'covariance matrix of Y\nwith standard deviation of {} is '
96               .format(sigma), np.round(get_original_and_noised_product_helper(
97                 sigma), 3))
98
99
100 # get_original_and_noised_product()
101
102
103 def run_all():
104     plot_all_eigenvals()
105     determine_original_eigenvals_amount()
106     get_original_and_noised_product()
107
108
109 # run_all()

```

3 SVM.py

```
1  import numpy as np
2  from matplotlib import pyplot as plt
3
4
5  BATCHES = [5, 50, 100]
6  ETAS = [1, 0.1, 10]
7  REG_PARAM = 0
8  ITERS = 150
9
10
11  image_size = 28
12  no_of_different_labels = 2
13  image_pixels = image_size * image_size
14  data_path = r'dataset\mldata'
15  train_data = np.loadtxt(data_path + r'\mnist_train.csv',
16                          delimiter=",")
17  test_data = np.loadtxt(data_path + r'\mnist_test.csv',
18                         delimiter=",")
19
20  fac = 0.99 / 255
21  all_train_imgs = np.asfarray(train_data[:, 1:]) * fac + 0.01
22  all_test_imgs = np.asfarray(test_data[:, 1:]) * fac + 0.01
23
24  all_train_labels = np.asfarray(train_data[:, :1])
25  all_test_labels = np.asfarray(test_data[:, :1])
26
27
28  class DigitsData:
29      """
30      Training and test sets for two digits to classify
31      """
32      def __init__(self, d_1, d_2):
33          binary_train_indices = np.logical_or((all_train_labels == d_1),
34                                              (all_train_labels == d_2))
35          binary_test_indices = np.logical_or((all_test_labels == d_1),
36                                             (all_test_labels == d_2))
37
38          binary_train_indices = binary_train_indices.reshape(
39              binary_train_indices.shape[0])
40          binary_test_indices = binary_test_indices.reshape(
41              binary_test_indices.shape[0])
42
43          x_train, y_train = all_train_imgs[binary_train_indices], \
44                          all_train_labels[binary_train_indices]
45          x_test, y_test = all_test_imgs[binary_test_indices], \
46                          all_test_labels[binary_test_indices]
47
48          y_train = y_train.reshape((len(y_train),))
49          y_test = y_test.reshape((len(y_test),))
50          y_train[y_train == d_1] = -1
51          y_test[y_test == d_1] = -1
52          y_train[y_train == d_2] = 1
53          y_test[y_test == d_2] = 1
54          self.x_train, self.y_train, self.x_test, self.y_test = \
55              x_train, y_train, x_test, y_test
56
57
58  def get_subgradient(w, reg_param, data, labels):
59      """
```

```

60     Gets a subgradient according the hinge loss
61     :param w: A d+1 by 1 numpy array where d is the dimension of each sample
62     i.e. the parameters of the classifier where the last element is the bias
63     :param reg_param: The regularization parameter in the sub-gradient
64     :param data: The samples to calculate the subgradient over
65     :param labels: The labels to calculate the subgradient over
66     :return: A subgradient according the hinge loss
67     """
68     sample = data[0]
69     sample_length = len(sample)
70     m = len(data)
71     only_w = w[: len(w) - 1]
72     b = w[-1]
73     sm = np.zeros(sample_length+1)
74     for i in range(m):
75         y_i = labels[i]
76         x_i = data[i]
77         extended_x_i = np.append(data[i], 1)
78         curr = y_i * extended_x_i
79         if y_i * (only_w @ x_i + b) < 1:
80             sm += curr
81     extended_w = np.append(only_w, 0)
82     return reg_param * extended_w - 1/m * sm
83
84
85
86 def gd(data, labels, iters, eta, w, reg_param):
87     """
88     The sub gradient descent algorithm
89     :param data: An n by d numpy array, where n is the amount of samples and
90     d is the dimension of each sample
91     :param labels: An n by 1 numpy array with the labels of each sample
92     :param iters: An integer that will define the amount of iterations
93     :param eta: A positive number that will define the learning rate
94     :param w: A d+1 by 1 numpy array where d is the dimension of each sample
95     i.e. the parameters of the classifier where the last element is the bias
96     :param reg_param: The regularization parameter in the sub-gradient
97     :return: A d+1 by 1 numpy array which contains the output of the sub
98     gradient descent algorithm over "iters" iterations and a list of the
99     vector result in each iteration
100    """
101    vecs = []
102    for iter in range(iters):
103        w += -eta * get_subgradient(w, reg_param, data, labels)
104        vecs.append(w.copy())
105    return 1/iters * np.sum(vecs, axis=0), vecs
106
107
108 def sgd(data, labels, iters, eta, w, batch):
109     """
110     The stochastic gradient descent algorithm
111     :param data: An n by d numpy array, where n is the amount of samples and
112     d is the dimension of each sample
113     :param labels: An n by 1 numpy array with the labels of each sample
114     :param iters: An integer that will define the amount of iterations
115     :param eta: A positive number that will define the learning rate
116     :param w: A d+1 by 1 numpy array where d is the dimension of each sample
117     i.e. the parameters of the classifier where the last element is the bias
118     :param batch: The amount of samples that the algorithm would draw and
119     use at each iteration
120     :return: A d+1 by 1 numpy array which contains the output of the
121     sub-gradient descent algorithm over "iters" iterations and a list of the
122     vector result in each iteration
123    """
124    vecs = []
125    for iter in range(iters):
126        samples_indices = np.random.choice(range(len(data)), batch)
127        samples = np.take(data, samples_indices, axis=0)

```

```

128     samples_labels = np.take(labels, samples_indices, axis=0)
129     v_t = get_subgradient(w, 0, samples, samples_labels)
130     w -= eta * v_t
131     vecs.append(w.copy())
132     return 1/iters * np.sum(vecs, axis=0), vecs
133
134
135 def get_errors(vecs, data, labels):
136     """
137     Gets a list of errors for each vector in vecs over the data and labels
138     :param vecs: The vecs to calculate their errors
139     :param data: The samples to calculate the errors over
140     :param labels: The labels to calculate the errors over
141     :return: A list of errors for each vector in vecs over the data and labels
142     """
143     errors = []
144     for vec in vecs:
145         errors.append(test_error(vec, data, labels))
146     return errors
147
148
149 def get_hypothesis(w):
150     """
151     Gets the classifying hypothesis of the given vector
152     :param w: The vector to get its classifying hypothesis
153     :return: The classifying hypothesis of the given vector
154     """
155     def h_w(x):
156         return np.sign(w[:-1] @ x + w[-1])
157     return h_w
158
159
160 def test_error(w, test_data, test_labels):
161     """
162     Returns a scalar with the respective 0-1 loss for the hypothesis
163     :param w: a d+1 by 1 numpy array where  $h_w(x) = \text{sign}(\langle w[:-1], x \rangle + w[-1])$ 
164     :param test_data: An n by d numpy array with n samples
165     :param test_labels: An n by 1 numpy array with the labels of the samples
166     :return: A scalar with the respective 0-1 loss for the hypothesis
167     """
168     h_w = get_hypothesis(w)
169     test_predictions = []
170     for sample in test_data:
171         test_predictions.append(h_w(sample))
172     result = test_labels - test_predictions
173     samples_amount = len(result)
174     non_zeros_amount = np.count_nonzero(result)
175     return non_zeros_amount / samples_amount
176
177
178 def plot_errors_helper(d_1, d_2, iters, gd_errors, sg_d_errors_lst, set_str, eta,
179                       reg_param):
180     iters_lst = range(iters)
181     plt.plot(iters_lst, gd_errors, label='GD')
182     for i in range(len(BATCHES)):
183         plt.plot(iters_lst, sg_d_errors_lst[i], label='SGD {}'.format(
184             BATCHES[i]))
185     plt.legend()
186     plt.xlabel('Iteration')
187     plt.ylabel('Loss')
188     plt.title('Losses of the GD and SGD algorithms\nclassifying the digits '
189             '{} and {} on the {} set\nper iteration, for {} '
190             'iterations\nwith learning rate of {}\nand regularization '
191             'parameter of {}'.format(d_1, d_2, set_str, iters, eta,
192                                     reg_param))
193     plt.show()
194
195

```

```

196 def plot_errors(d_1, d_2, train_samples, train_labels,
197                 test_samples, test_labels, iters, eta, w, reg_param):
198     w_gd, vecs_gd = gd(train_samples, train_labels, iters, eta, w, reg_param)
199     gd_train_errors = get_errors(vecs_gd, train_samples, train_labels)
200     gd_test_errors = get_errors(vecs_gd, test_samples, test_labels)
201     sgd_train_errors_lst = []
202     sgd_test_errors_lst = []
203     for batch in BATCHES:
204         w_sgd, vecs_sgd = sgd(train_samples, train_labels, iters, eta, w, batch)
205         sgd_train_errors_lst.append(get_errors(vecs_sgd, train_samples,
206                                                train_labels))
207         sgd_test_errors_lst.append(get_errors(vecs_sgd, test_samples,
208                                                test_labels))
209     plot_errors_helper(d_1, d_2, iters, gd_train_errors, sgd_train_errors_lst,
210                       'training', eta, reg_param)
211     plot_errors_helper(d_1, d_2, iters, gd_test_errors, sgd_test_errors_lst,
212                       'test', eta, reg_param)
213
214
215 def classify_digits(d_1, d_2):
216     """
217     Run all items for the given digits
218     :param d_1: First digit
219     :param d_2: Second digit
220     :return: None
221     """
222     sets = DigitsData(d_1, d_2)
223     sample = sets.x_train[0]
224     w_length = len(sample) + 1
225     zero_w = np.zeros(w_length)
226     for eta in ETAS:
227         plot_errors(d_1, d_2, sets.x_train, sets.y_train,
228                   sets.x_test, sets.y_test, ITERS, eta, zero_w, REG_PARAM)

```


4 ex6 runme.py

```
1  from SVM import classify_digits
2
3
4  def main():
5      classify_digits(0, 1)
6      classify_digits(3, 5)
7
8
9  if __name__ == "__main__":
10     main()
```