

# תרגיל בית מספר 3

מבוא ללמידה חישובית

בנימין אסתרליס ותומר סמרה

315636761 321164923

tomersamara@mail.tau.ac.il esterlis@mail.tau.ac.il

הכתובת של קבצי הקוד + התמונות בצבע: specific/a/home/cc/students/math/esterlis/Downloads/ML-HW3

# 1 חלק תיאורטי

1. ראשית, ניזכר בעובדה המכיתה כי ה  $VCdim$  של נירון בודד הינו  $d+1$  לפי *Radons' theorem*. נרצה להראות שה  $VCdim$  של רשת נירונים היא לפחות ה  $VCdim$  של נירון בודד: לשם כך נגדיר את רשת הנירונים הבאה, אשר תתנהג בדיוק כמו נירון בודד: את הערך של הנירון הראשון בשכבה הראשונה נעביר האלה לכל שאר הנירונים ע"י מכפלה של השכבה הראשונה במטריצת היחידה, והפעלת  $sign$  קורדינאטה-קורדינאטה. נעשה זאת על כל השכבות מהשנייה ועד השכבה ה  $L - 1$ . לבסוף בשכבה ה  $L$  נעביר כפלט את הסימן של הנירון הראשון בשכבה ה  $L - 1$ , ולכן קיבלנו כי הנירון הבודד בשכבה האחרונה מחזיר לנו את הסימן של הנירון הראשון בשכבה הראשונה.

הראנו שהרשת שלנו מתנהגת בדיוק כמו נירון בודד ולכן כל מה שניתן לסווג ע"י נירון בודד ניתן לסווג ע"י הרשת שבנינו, ולכן ה  $VCdim$  של רשתות נירונים  $\leq VCdim$  של נירון בודד.

2. (א) מכון שהפונקציית המטרה היא פונקציה מ  $\mathcal{R}^d \rightarrow \mathcal{R}^d$  אזי מרחב ההשערות שלנו הוא כל הפונקציות מהסוג הזה, כעת לפי הגדרת  $\mathcal{H}$  אנו מקבלים פונקציות מ  $\mathcal{R}^d \rightarrow \mathcal{R}$  ולכן בעצם אנו צריכים לשרשר את  $\mathcal{H}$ ,  $d$  פעמים ולכן נקבל כי מחלקת  $\mathcal{H}^d$  ההיפותזות שלנו היא:

נחסום כעת את קצב הגדילה, לפי *Lemma 7.2* והחישוב שעשינו בתרגול מתקבל כי קצב הגדילה של נירון בודד הוא:  $\left(\frac{em}{d+1}\right)^{d+1} \leq$  ולכן לפי *Lemma 7.3* מתקבל כי החסם הוא:

$$\pi_{\mathcal{H}^d}(m) \leq (\pi_{\mathcal{H}}(m))^d \leq \left(\left(\frac{em}{d+1}\right)^{d+1}\right)^d = \left(\frac{em}{d+1}\right)^{d \cdot (d+1)}$$

(ב) מכון ש  $\mathcal{C}$  הינה מחלקה של כל רשתות הנירונים עם  $L$  שכבות ו  $d$  נירונים, ובסעיף הקודם הראנו כי כל שכבה הינה מהצורה הבאה:  $\mathcal{H}^d$

ולכן מתקבל כי המחלקה הנ"ל היא:  $\mathcal{H}^{d \cdot \mathcal{H}^d}$ , כאשר בחזקה יש  $\mathcal{L} - 1$  פעמים  $\mathcal{H}^d$ . נחסום כעת את קצב הגדילה, לפי הנימוקים של התרגיל הקודם מתקבל כי קצב הגדילה של נירון בודד הוא:  $\left(\frac{em}{d+1}\right)^{d+1} \leq$  ולכן לפי *Lemma 7.4* מתקבל כי החסם הוא:

$$\pi_{\mathcal{C}}(m) \leq \pi_{layer1}(m) \cdot \dots \cdot \pi_{layer\mathcal{L}-1}(m) \cdot \pi_{layer\mathcal{L}}(m) \leq \left(\frac{em}{d+1}\right)^{d(d+1)(\mathcal{L}-1)+(d+1)}$$

(ג) ניתן לראות כי יש לנו  $\mathcal{L} - 1$  שכבות בעלות  $d$  נירונים בכל שכבה, בנוסף לכך יש לנו בשכבה האחרונה רק נירון אחד, ובכל נירון יש לנו  $d + 1$  וולכן בסה"כ נקבל כי מספר הפרמטרים הדרוש הוא:

$$N = (d + 1) \cdot ((\mathcal{L} - 1) \cdot d + 1) = d \cdot (d + 1) \cdot (\mathcal{L} - 1) + (d + 1)$$

(ד) בנוס, מחוסר זמן לצערי לא היה זמן לעשות

(ה) בעקבות ההוכחת שהראנו עד הלום מתקיים:

$$2^{VCdim(\mathcal{C})} = \pi_{\mathcal{C}}(m) \leq \left(\frac{em}{d+1}\right)^{d \cdot (d+1) \cdot (\mathcal{L}-1) + (d+1)} = \left(\frac{em}{d+1}\right)^N \leq (em)^N$$

וכאשר נציב  $VCdim(\mathcal{C}) = m$ , ובעקבות מה שהראתי החלק הקודם של שאלה וסעיף  $d$  מתקבל כי  $m \leq 2N \log_2(en)$  ולכן בסה"כ  $VCdim(\mathcal{C}) \leq 2N \log_2(en)$ .

3. (א) האלגוריתם:

האלגוריתם  $SGD$  החדש פועל באותו האופן של האלגוריתם הרגיל, רק בשלב של בחירת הנקודה נבחר את הנקודה באופן הבא:

$$y_{t+1} = x_t - \nabla f(x_t)$$

$$x_{t+1} = \Pi_{\mathcal{K}}(y_{t+1})$$

ואת  $\Pi_{\mathcal{K}}(x)$  באופן הבא:

$$\Pi_{\mathcal{K}}(x) = \begin{cases} x & x \in \mathcal{K} \\ \frac{x}{\|x\|} \cdot R & \text{otherwise} \end{cases}$$

שזה בעצם לבחור את הנקודה הכי קרוב לנקודה  $x$  אשר בקבוצה הקמורה  $\mathcal{K}$ , המעגל ברדיוס  $R$ .

(ב) בונס\_בני

(ג) ראשית ניתן לראות כי בהוכחה שעשינו בכיתה עד שיווין (12) לא מתעסקים באלגוריתם  $SGD$  עצמו ולכן החלק הזה בהוכחה נשאר גם פה ומתקבל כי:

$$\mathbb{E}[(\bar{W}) - f(w^*)] \leq \frac{1}{T} \sum_i \mathbb{E}[V_t(W_t - w^*)]$$

מכיוון שאנו משתמשים ב  $SGD$  עם  $projection$  מהדרך שבה אנו מקדמים את  $W_t$  מתקבל ש  $W_{t+1} = \Pi_{\mathcal{K}}(W_t - \eta V_t)$  ולכן מתקיים  $\|W_{t+1} - w^*\| = \|\Pi_{\mathcal{K}}(W_t - \eta V_t) - w^*\|$  לפי סעיף ב' ומכיוון ש  $\mathcal{K}$  הינה קמורה מתקיים

$$\begin{aligned} \|W_{t+1} - w^*\|_2^2 &\leq \|W_t - \eta V_t - w^*\|_2^2 = \\ &= \|W_t - w^*\|_2^2 + \eta^2 \|V_t\|_2^2 - 2\eta (W_t - w^*)^T V_t \end{aligned}$$

והעברת אגפים תיתן לנו :

$$(W_t - w^*)^T V_t \leq \frac{\|W_t - w^*\|_2^2 - \|W_{t+1} - w^*\|_2^2}{2\eta} + \frac{\eta^2 \|V_t\|_2^2}{2}$$

וכעת המשיך ההוכחה הוא זהה לחלוטין להוכחה שראינו בכית החל מסעיף 1

**Bonus a Deserve I**

stabbing without week the through it Making For

fork a with someone

(ד) נראה כי הטענה אינה נכונה בעזרת דוגמא נגדית:

נבחר  $(x, z) := K_2(x, z) \forall x, z$  ואת  $K_1$  באופן שרירותי ולכן נקבל כי  $K(x, z) = \frac{K_1(x, z)}{0} \forall x, z$  ומכיון שחלוקה ב0 איננה מוגדרת אזי גם  $K$  אינו מוגדר ולכן בפרט  $K$  אינו גרעיון מוגדר חיובית.

(ה) ידוע לנו ש  $K = aK_1$  ולכן מכיוון ש  $K_1$  מוגדרת חיובית לכל וקטור  $v$  מתקיים  $v^T \cdot K_1 \cdot v > 0$  ולכן מתקיים  $v^T \cdot K \cdot v = v^T \cdot aK_1 \cdot v > 0$  קיבלנו ש  $K$  מוגדרת חיובית.

4. (א)

(ב) ראשית נסמן:

$$j^* = \underset{j \in [K]}{\operatorname{argmax}} (w_j \cdot x_i - w_{y_i} \cdot x_i + \chi_{(y_i \neq j)})$$

כעת נגזור את  $l$ :

על מנת לחשב הדרוש נפצל למקרים הבאים:

$$j \neq j^* \wedge j \neq y_i \quad \text{i.}$$

מכיון ש  $l$  אינה תלויה כלל ב  $j$  נקבל כי הנגזרת היא 0 במקרה זה.

ii.  $j \neq j^* \wedge j = y_i$   
מכוון ש  $l$  תלוי ב  $j$  נקבל

$$l(w_1, \dots, w_K, x_i, y_i) = w_{j^*} \cdot x_i - w_j \cdot x_i + \chi_{j \neq j^*}$$

ולכן גזירה לפי  $j$  נקבל כי הנגזרת היא:  $-x_i$ .

iii.  $j = j^* \wedge j \neq y_i$   
מכוון ש  $l$  תלוי ב  $j$  נקבל כי:

$$l(w_1, \dots, w_K, x_i, y_i) = w_j \cdot x_i - w_{y_i} \cdot x_i + \chi_{(j \neq y_i)} = w_j \cdot x_i - w_{y_i} \cdot x_i + 1$$

ולכן גזירה לפי  $j$  נקבל כי הנגזרת היא:  $x_i$ .

iv.  $j = j^* \wedge j = y_i$   
מכוון ש  $l$  תלוי ב  $j$  נקבל כי:

$$l(w_1, \dots, w_K, x_i, y_i) = w_j \cdot x_i - w_j \cdot x_i + \chi_{(j \neq y_i)} = w_j \cdot x_i - w_j \cdot x_i + 1$$

ולכן גזירה לפי  $j$  נקבל כי הנגזרת היא:  $x_i - x_i = 0$ .

ולכן בסה"כ:

$$\nabla_{w_j} l = \begin{cases} -x_i & j \neq j^* \wedge j = y_i \\ x_i & j = j^* \wedge j \neq y_i \\ 0 & otherwise \end{cases}$$

ולכן ניתן לחשב את ה  $sub$  descent gradient באופן הבא:

$$\nabla_{w_j} f = w_j + \nabla_{w_j} l = \begin{cases} w_j - C' x_i & j \neq j^* \wedge j = y_i \\ w_j + C' x_i & j = j^* \wedge j \neq y_i \\ w_j & otherwise \end{cases}$$

כאשר  $C' = \frac{C}{m}$ .

וכעת נתאר האלגוריתם  $Stochastic$  subgradient descent כאשר  $C, T, \eta, \{x_i\}_{i=1}^m, \{y_i\}_{i=1}^m$  נתונים לנו:

i. נאתחל את  $w_{i0} = 0 \quad \forall 1 \leq i \leq k$

ii. לכל  $t \in [T]$  נבצע את הבאים:

א'. נבחר באופן אקראי עם התפלגות אחידה  $i \in [m]$

ב'. נחשב  $j^*$  באופן הבא:

$$\operatorname{argmax}_{j \in [k]} (w_j \cdot x_i - w_{y_i} \cdot x_i + \chi_{(j \neq y_i)})$$

ג'. לכל  $j \in [k]$  נבצע את הבאים:

$$d'. \quad w_{j,t+1} = (1 - \eta) \cdot w_{j,t}$$

ה'. אם  $j^* \neq y_i$

$$v'. \quad w_{j^*,t+1} = w_{j^*,t} - \eta \cdot C \cdot x_i$$

$$z'. \quad w_{j^*,t+1} = w_{j^*,t} + \eta \cdot C \cdot x_i$$

iii. ונחזיר  $(w_{(1,T)}, \dots, w_{(k,T)})$ .

והסייוג של  $x$  ל  $\{1, \dots, k\}$  יתצבע ע"י החישוב:

$$\operatorname{argmax}_{j \in [k]} \langle w_j, x \rangle$$

5. על מנת שנשנה את האלגוריתם הנוכחי שיתשמש ב  $Kernel$  נשים לב כי בכל שלב באלגוריתם  $stochastic subgradient descent$

$$w_{jt} = \sum_{i=1}^t \beta_{ji} \cdot x_i$$

ה  $w_{j,t}$  הינו צריוף לינארי של הדגימות שהיו קודם לכן ולכן מקבל כי:  $\beta_{ji} = (1 - \eta)^{t-1} \cdot A_{ji} \eta \cdot C$  כאשר  $\alpha_{ij} = \begin{cases} 1 & j^* \neq y_i \wedge j = y_i \\ -1 & j^* \neq y_i \wedge j = j^* \\ 0 & otherwise \end{cases}$  ו  $j^*$  הינו התקבל בשלב ה  $i$ .  
 מכוון ש  $\beta_{ji}$  משתנה במהלך הריצה ולפי מה שראינו בתרגיל הקודם עבור  $Kernel Perceptron$  ולכן מתקיים כי:

$$\langle w_{jt}, x \rangle = \left\langle \sum_{i=1}^t \beta_{ji} \cdot x_i, x \right\rangle = \sum_{i=1}^t \beta_{ji} \langle x_i, x \rangle$$

והאלגוריתם SGD אשר משתמש ב  $Kernel$ :

(א) נגדיר מטריצה חדשה  $A_{k \times m}$  להיות מטריצית אפסים שתייצג את המקדמים  $\alpha_{ji}$

i. נחשב את  $j^*$  באופן הבא:

$$j^* = \underset{j \in [k]}{\operatorname{argmax}} \left( \sum_{i=1}^{t-1} [\beta_{ji} \cdot K(x_i, x_t) - \beta_{y_t i} \cdot K(x_i, x_t)] + \chi_{(j \neq y_t)} \right)$$

ii. כעת אם  $j^* \neq y_t$  נציב: (בתאים האחרים נשאר 0)

$$A_{j^* t} = -1$$

$$A_{y_t, t} = 1$$

$$\forall \begin{cases} 1 \leq i \leq m \\ 1 \leq j \leq k \end{cases} \quad \beta_{ji} = (1 - \eta)^{t-1} \cdot A_{ji} \eta \cdot C$$

ג'. כעת נעדכן את  $\beta_{ji}$

$$\text{iii. כעת נגדיר מטריצה חדשה } B_{k \times m} \text{ כאשר } B_{ji} = \beta_{ji} \quad \forall \begin{cases} 1 \leq i \leq m \\ 1 \leq j \leq k \end{cases} \text{ כאשר } \beta_{ji} \text{ הינו לאחר העדכון.}$$

$$\underset{j \in [k]}{\operatorname{argmax}} \left( \sum_{i=1}^m \beta_{ji} \cdot K(x_i, x) \right)$$

לבסוף סיווג הנקודה החדשה  $x$  יתבצע ע"י החישוב:

6. נוכיח את הטענה באינדוקציה על  $d$

(א) צעד בסיס ( $d=1$ ):

צריך כי  $h : \{0, 1\} \rightarrow \{0, 1\}$  היה עץ החלטות מגובה 2 ולכן נגדיר אותו באופן הבא:  
 עץ החלטות שלנו יהיה בעל שורש יחיד מצורה:  $[(x_1 = 0)?]$  ולכן זה עץ מגובה 2, והוא מקיים את התנאים הדרושים

(ב) צעד אינדוקציה ( $d = 1$ ):

נניח כי מתקיים עבור  $d = n$  ונוכיח עבור  $n + 1$  ולכן נגדיר את עץ ההחלטות  $T$  באופן הבא:

i. נבנה את העץ מגודל  $n + 1$  עבור הקורדינאטות  $x_1, \dots, x_n$ , מכוון שיש לנו  $n$  קורדינאטות מתקבל כי תנאי ההנחה מתקיימים וניתן לבנות עץ החלטות  $T'$  מגודל  $n + 1$  כשבשכבה ה  $i$  נשאל את השאלה  $x_i = 0$ ? ונתקדם לפיה.

ii. כעת נוסיף לעץ  $T'$  שכבה חדשה, לכל עלה נוסיף צומת חדשה על מנת לסווג גם לפי הקורדינאטה האחרונה,  $x_{n+1}$ , כך שצומת תהיה מהצורה:  $(x_{n+1} = 0)$ ?

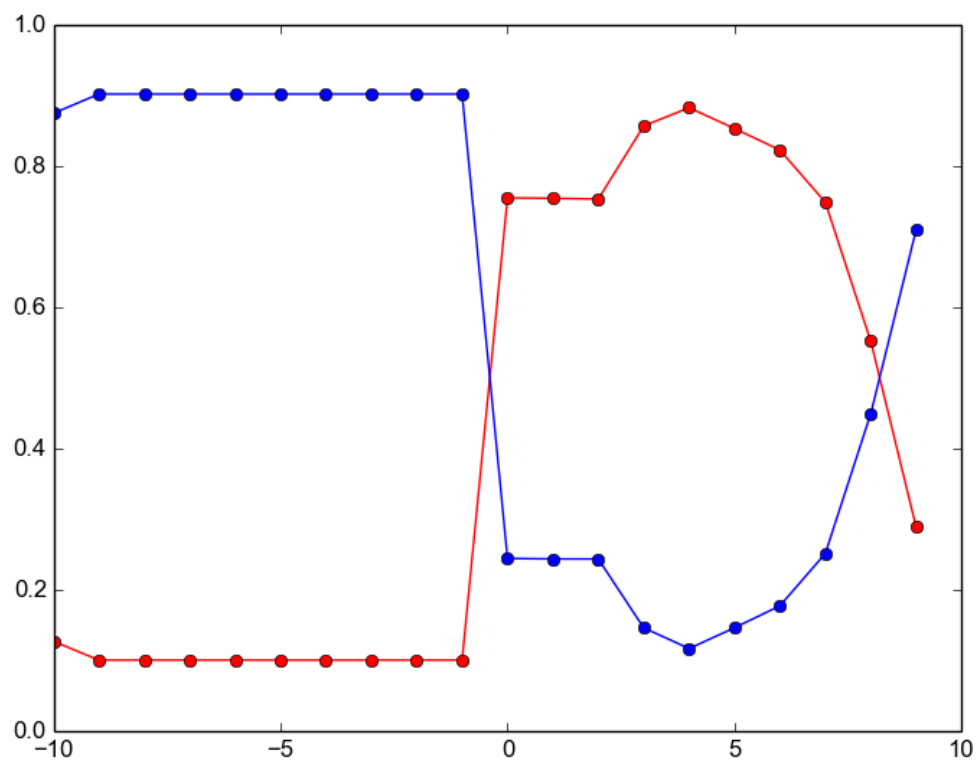
(ג) כעת קיבלנו את העץ  $T$  אשר הוא מסווג מהצורה:  $h : \{0, 1\}^{n+1} \rightarrow \{0, 1\}$ , בעקבות בניית  $T$  מתקבל כי העץ הוא מגול כעת צריך לנתץ את הקבוצה  $\{0, 1\}^d = S$  ונקבל כי  $VCdim \leq 2^d$ . כדורוש.  $n + 1 + 1 = n + 2$ .

ולכן עבור הקבוצה  $S$  הנ"ל נראה כי לכל  $v \in \{0, 1\}^{2^d}$ , סיווג כלשהו, קיימת פונציקה  $h \in \mathcal{H}$  כך ש  $h(s_i)_{s_i \in S} = v_i$  כאשר  $\mathcal{H}$  הינה מחלקת עצי ההחלטה.

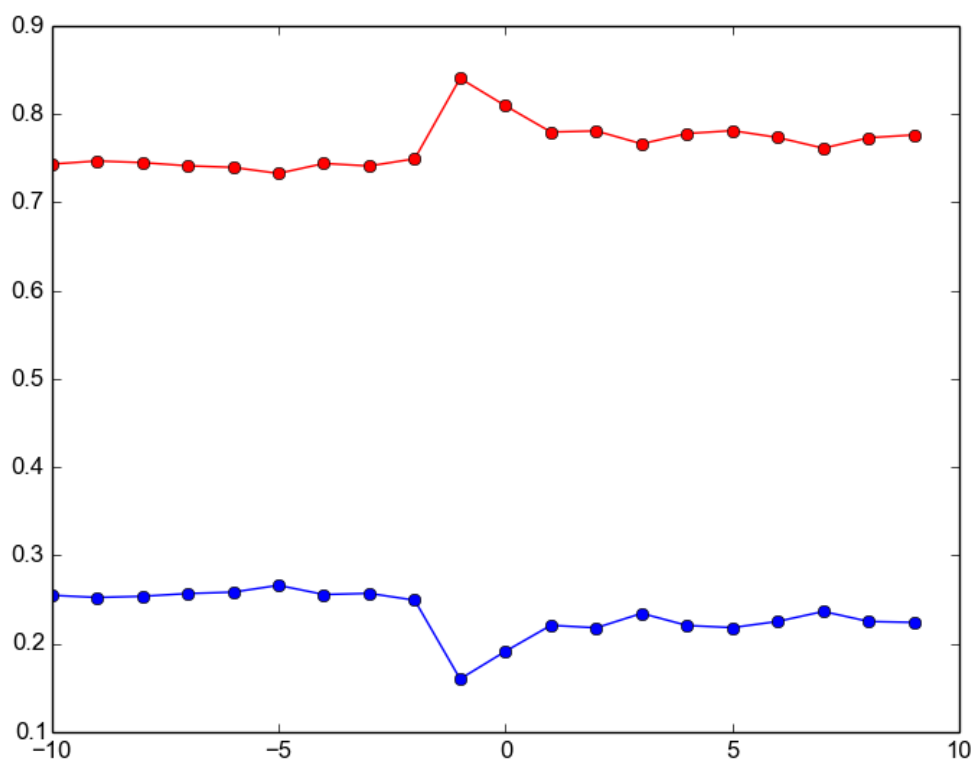
יהי  $v \in \{0, 1\}^{2^d}$  ונבנה לו את הפונציקה  $f$ , שתשרה את הסיווג, באופן הבא:  $h(s_i) = (v)_i$  כאשר  $h : \{0, 1\}^{2^d} \rightarrow \{0, 1\}$ . ניתן לבנות פונקציה כזו מכוון שבחלק הקודם של השאלה הוכחנו כי קיים עץ אשר משרה את הפונקציה הזו, ולכן הראנו כי זה מתקיים לכל  $v \in \{0, 1\}^{2^d}$  ו  $\Pi_{\mathcal{K}}(S) = 2^{2^d}$  ולכן מתקבל כי  $VCdim \leq 2^d$ .

## 2 חלק מעשי

1. (א) ניתן לראות על פי התוצאות שה $\eta$  הטובה ביותר מתקבלת ב $10^{-6}$ .  
 ה $C$  האופטימלי מתקבל ב0.1.  
 והדיוק הוא:0.9075

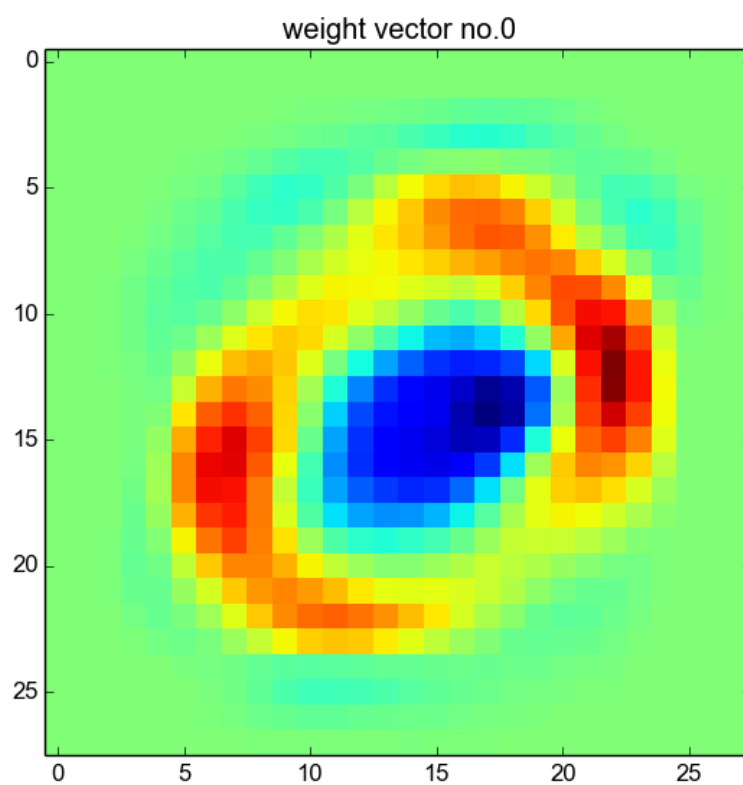


i.



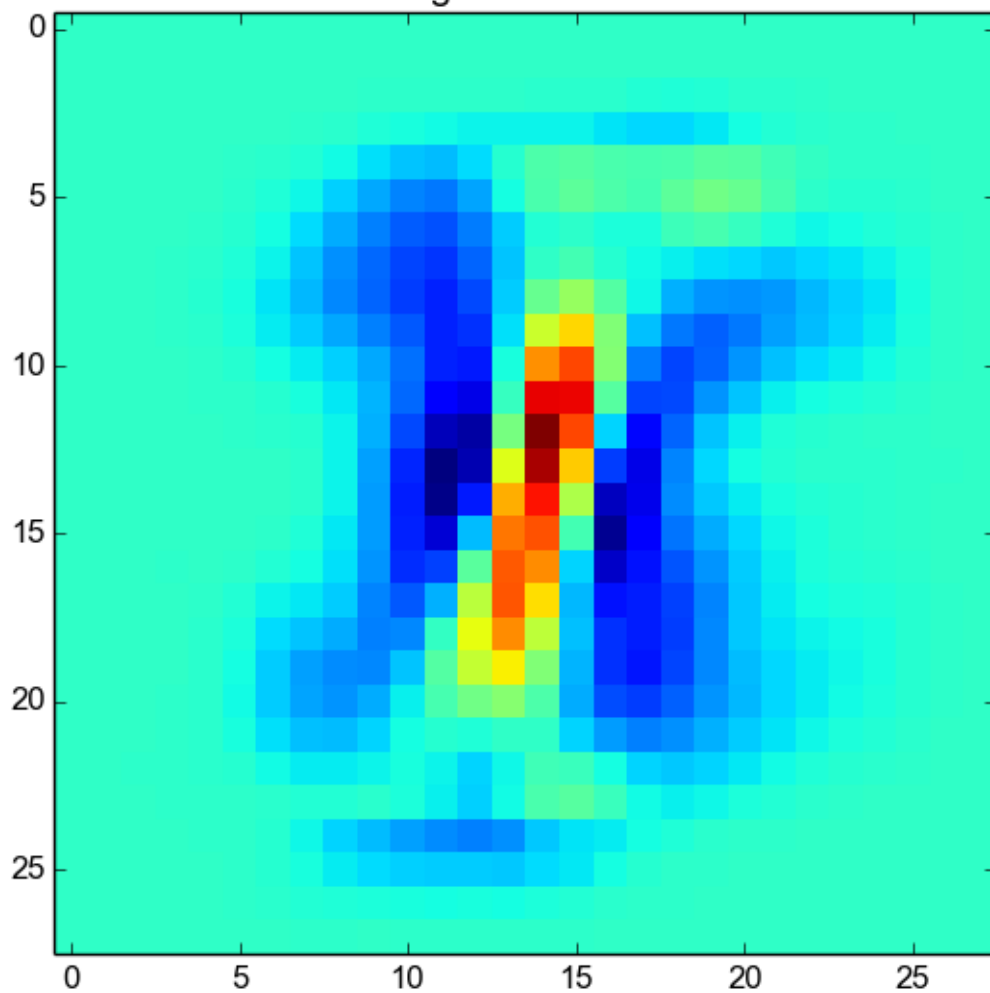
ii.

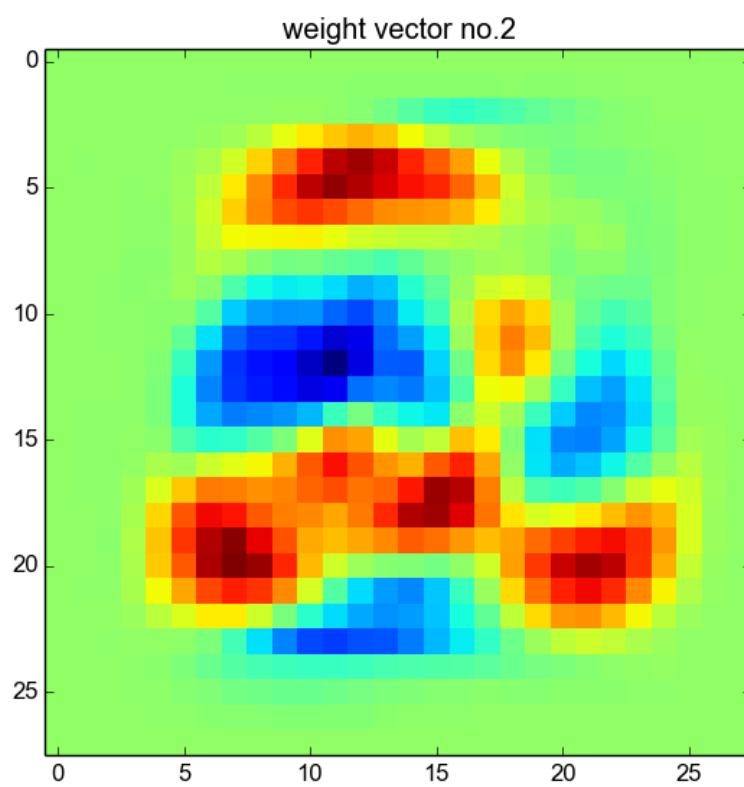




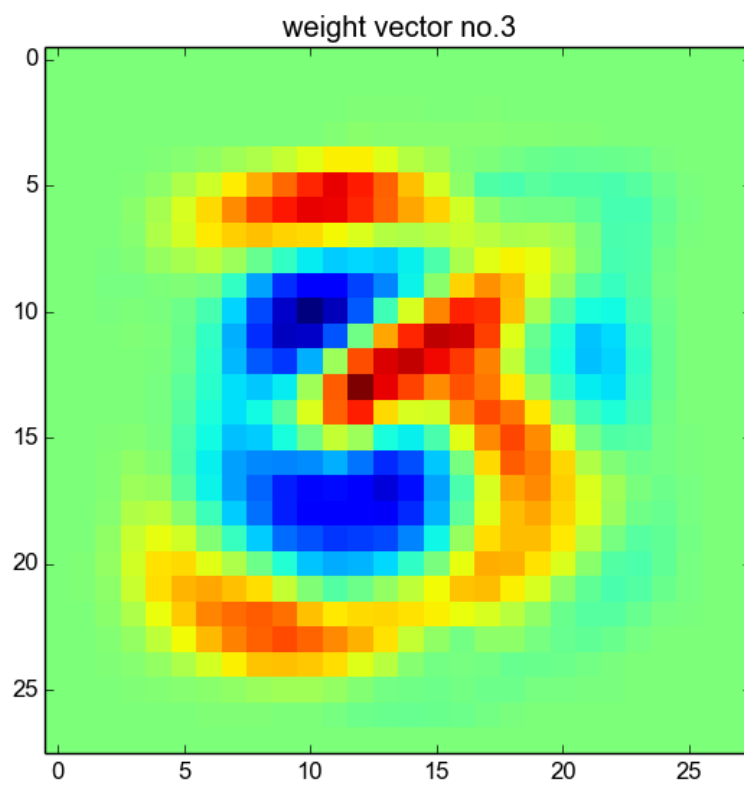
i. (2)

weight vector no.1

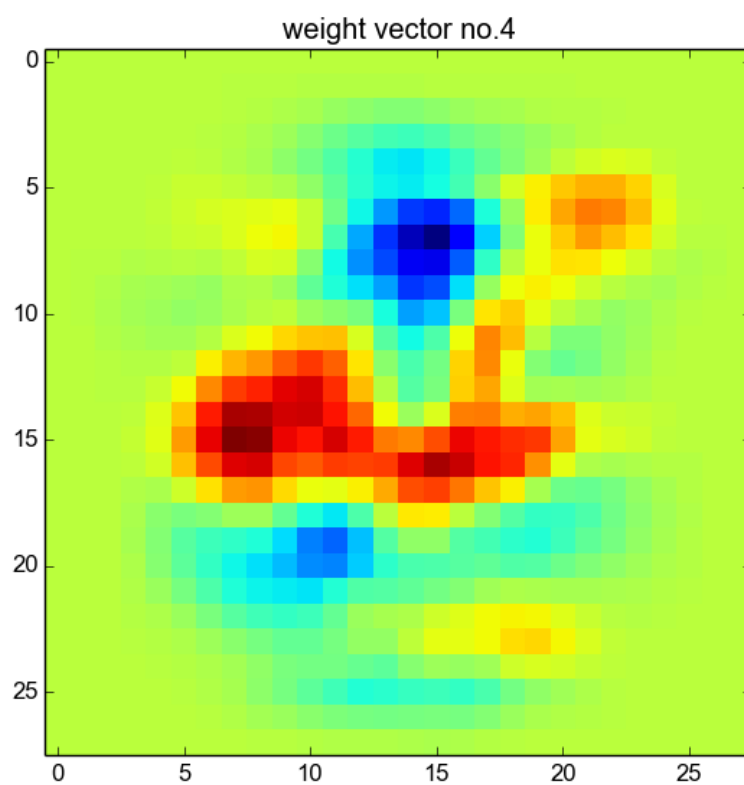




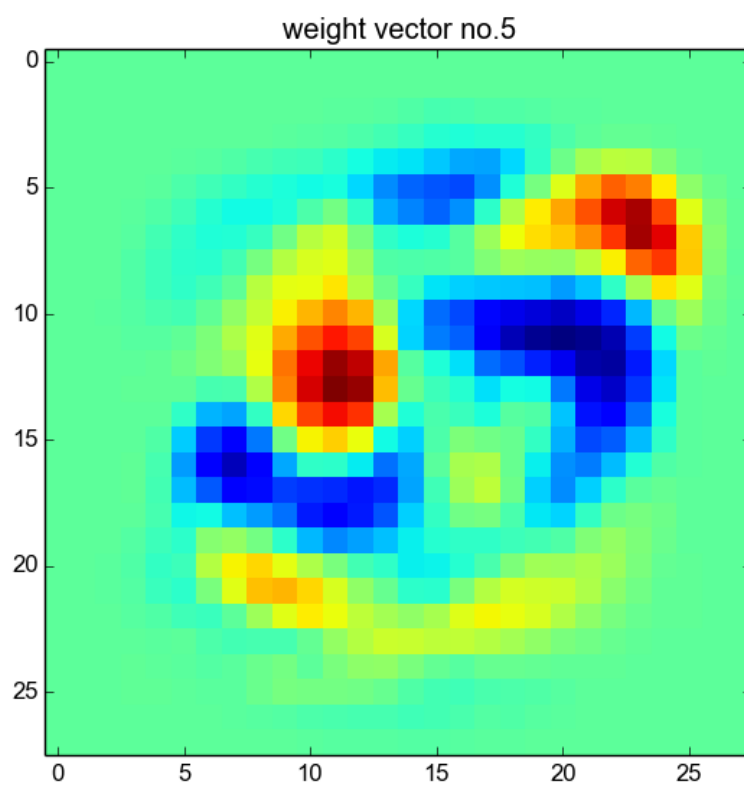
iii.

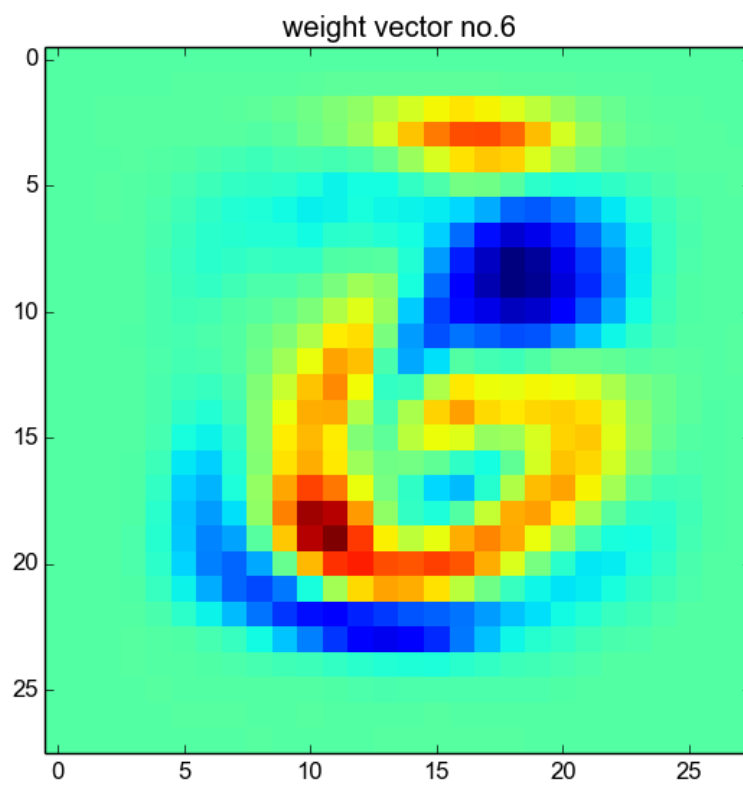


iv.

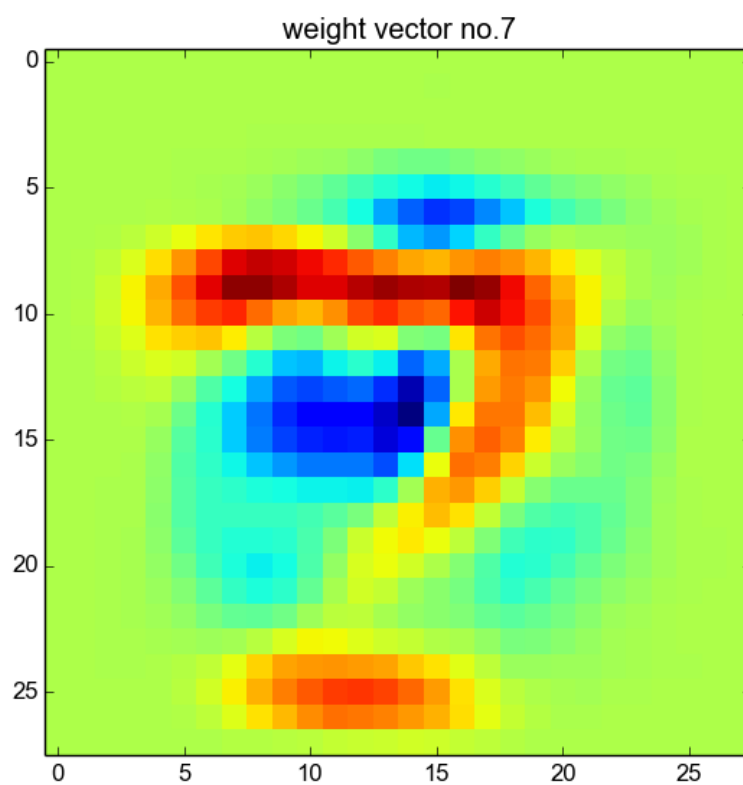


v.



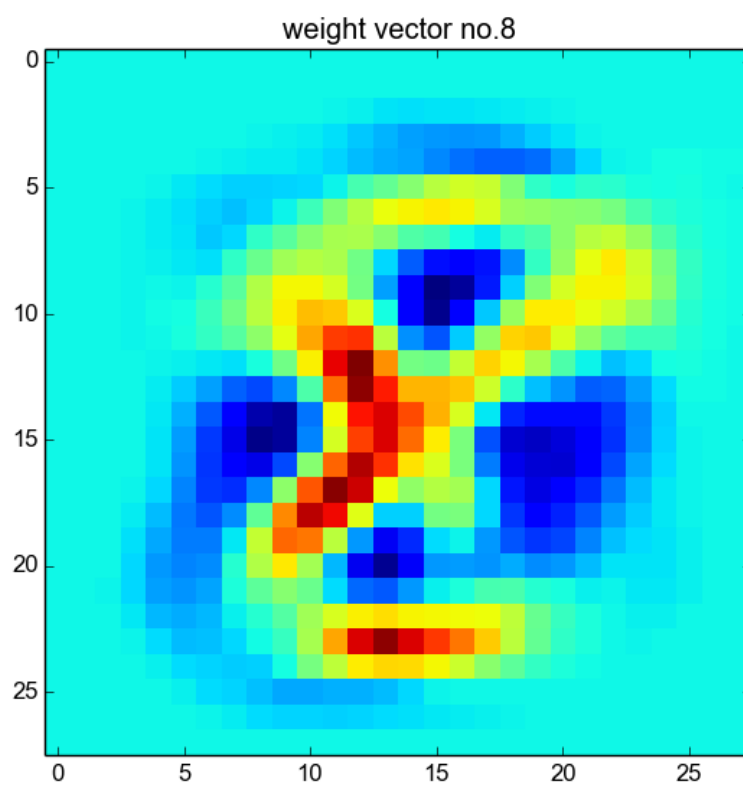


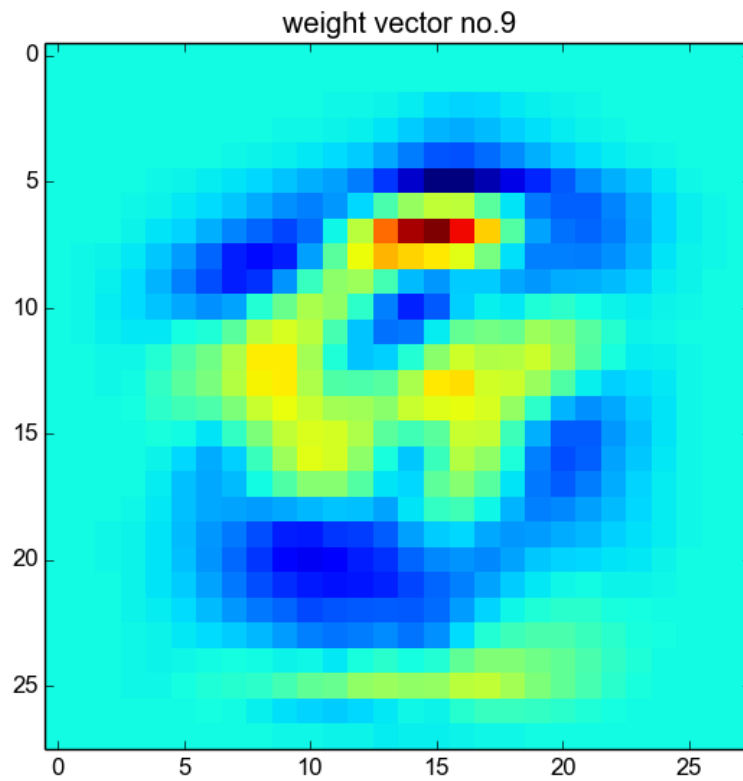
vii.



viii.







x.

(ג) הדיוק שהתקבל לנו על *test set* הינו 0.9167

2. (א)  $\eta$  הטוב ביותר התקבל ב  $10^{-3}$  עבור  $T = 500$  ובהינתן  $\eta$  זה מתקבל כי  $C = 10^9$ .  
 הדיוק שהתקבל עבור  $\eta$  זו הינו: 0.8224 ועבור  $C$  זה הינו: 0.998 על הדגימה ועל ה *training set* הוא: 0.8407.

(ב) הדיוק על *test set* הוא: 0.7945.