

Neural Network For Images

29 במרץ 2023

Ex1

איתמר שכטר

315092759

itamar.sh

שאלות פרקטיות:

שאלה פרקטית ראשונה *Architecture*:

בחרתי להשתמש בארכיטקטורה בסיסית ולהעלות את מספר *channels* בכל השכבות.

הארכיטקטורה הבסיסית הייתה בנויה מ:

(1) שכבת קונבולוציה כשגודל הפילטר הוא 5 על 5 ובעומק 3 כי אנחנו קיבלנו קלט של

RGB. עם אקטיביציית *relu*.

למעשה בשלב הזה אנחנו לוקחים תמונה בגודל 32 על 32 על 3 והופכים אותה ל-28 על 28

על מספר הערוצים שנבחר. (*padding=valid*)

(2) שכבת *max pooling* עם *stride* 2 על 2 ובגודל 2 על 2.

בשלב זה אנחנו מצמצמים שוב את מימדי התמונה ל-14 ל-14 ומשאירים את כמות

הערוצים אותו דבר.

(3) שוב שכבת קונבולוציה כשגודל הפילטר הוא 5 על 5 ומספר הערוצים ככמות הערוצים

שהיו בפלט השכבה הראשונה. עם אקטיביציית *relu*.

מימדי התמונה ישתנו ל-10 על 10.

(4) שוב שכבת *max pooling* עם *stride* 2 על 2 ובגודל 2 על 2.

בשלב זה אנחנו מצפים לקבל פלט בגודל 5 על 5 כפול מספר הערוצים בפלט השכבה

הקודמת.

(5) בשלב זה משטחים את התמונה לוקטור ומכניסים את הקלט לשכבת *FC* עם גודל פלט

10 כגודל כמות המחלקות שלנו שלהם אנחנו רוצים לבנות וקטור הסתברות. עם אקטיביציית

relu.

- אני מעיר שיש לנו הייפר פרמטר של קצה הלמידה $learning\ rate = 0.001$, מספר

האפוקים שבחרתי הוא 15.

כל המשחק מעתה והלאה יהיה בכמות הערוצים שנבחר ל-25 שכבות הקונבולוציה.

בכוונה אנחנו שומרים על פתרון יחסית פשוט, בלי הרבה שכבות קונבולוציה קטנות בגודל

3 על 3 כי חשבתי שהרבה שכבות קטנות יהיה קשה לבחור היפר פרמטרים שונים ולהסיק

מסקנות בצורה ברורה. למרות שהבנתי שרוב הרשתות הקיימות בסוף מתכנסות לשכבות 3

על 3 כי זה נותן יותר *perceptive field* בפחות כוח עיבוד.

אני לקחתי מערכת עם מעט מאוד פרמטרים.

על כמות הפרמטרים אני שולט בכמות הערוצים בשכבות הקונבולוציה.

בריצה הראשונה התחלתי עם:

3 ערוצי פלט בשכבת קונבולוציה הראשונה.

8 ערוצי פלט בשכבת הקונבולוציה השנייה.

נחשב את כמות הפרמטרים:

$conv1$: גודל כל קרנל הוא 5 על 5 על 3 ולכן גודל קרנל הוא 75 ונוסיף לזה באייס אז

76. יש לנו 3 ערוצים פלט ולכן $228 = 3 * 76$.

$conv2$: גודל כל קרנל הוא 5 על 5 על 3 ולכן שוב נקבל 76 לקרנל יחיד. יש לנו 8 ערוצי

פלט ולכן 8 קרנלים, אז החישוב הוא: $608 = 8 * 76$.

FC : נבחין שגודל הקלט הוא גודל התמונה שזה 5 על 5 בשלב הזה (ראו סעיף 4 למעלה)

והעומק הוא לפי ערוצי הפלט של השכבה השנייה שזה 8.

ולכן נקבל שיש לנו סה"כ וקטור קלט באורך $200 = 8 * 5 * 5$. הפלט שלנו הוא בגודל

10 ולכן כמות הפרמטרים היא מטריצה בגודל $200 * 10$ ועוד באייס באורך הפלט שהוא 10

ולכן סה"כ $2010 = 10 * 201$.

סה"כ: $2846 = 2010 + 608 + 228$.

בכל ריצה מהריצות הבאות הכפלתי ב2 את כמות הערוצים מהריצה שלפניה בקונבולוציה,

בעיקר כדי להגיע יחסית מהר ל $over fit$.

מדדתי בכל 2000 מיני באצ'ים את ה $train loss$, זה פחות מדויק כי יש פה רנדומליות

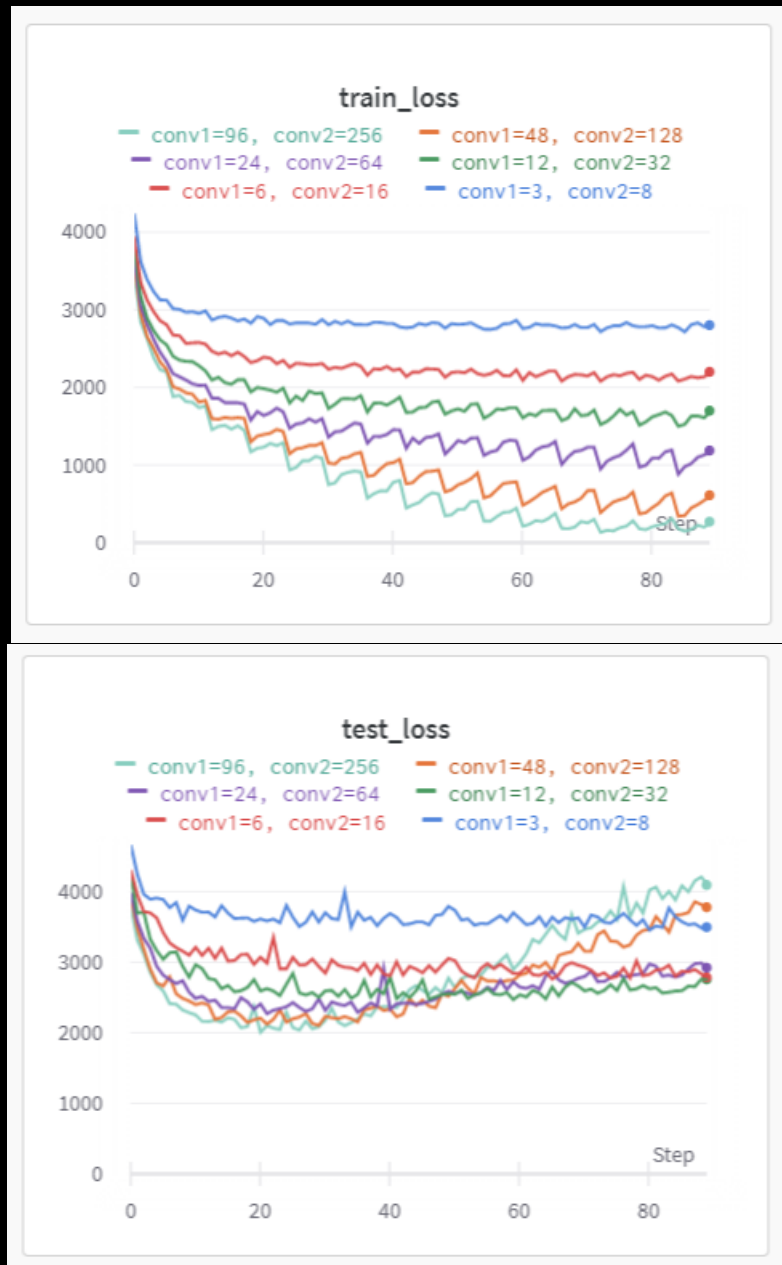
אבל בסה"כ אנחנו מצפים לראות גרף של ירידה.

אם לא נראה גרף של ירידה ב $train$ סימן שיש לנו $under fit$.

בנוסף מדדתי את ה $test loss$ על פני כל ה $test data$ בכל 2000 מיני באצ'ים כדי לראות

ירידה גם פה. במידה ונתחיל לראות נקודה בה ה $test$ מתחיל לעלות נדע שהגענו ל $over fit$.

אלה התוצאות שקיבלתי:



אני אסביר:

כל צבע כאן מציג אחת מ-6 הרשתות שאימנתי, השם שלהן הוא לפי כמות הצ'אנלים שיש בפלט שכבת הקונבולוציה הראשונה ופלט שכבת הקונבולוציה השנייה. (יש רק שני שכבות) לדוגמא $conv1 = 24, conv2 = 64$ מציג רשת עם שכבת קונבולוציה ראשונה בעלת 24 קרנלים ופולטת קלט עם 64 ערוצים. כלומר אם הקלט הוא תמונות בגודל (32, 32, 3) אז הפלט הוא טנזורים בגודל: (28, 28, 24). כי כל השכבות הן עם קרנל 5 על 5. והשכבה

השנייה (אחרי $maxpool$) תקבל קלט (14, 14, 24) ותחזיר (10, 10, 64).

בגרף הלמידה ניתן לראות שהצבע הכחול ($conv1 = 3, conv2 = 8$) בעל הכי פחות פרמטרים ללמידה הוא כמעט ולא לומד ואכן גם ה- $test loss$ שלו די נשאר במקום. הגרף האדום, ($conv1 = 6, conv2 = 16$), עם כפול ערוצים בכל שכבת קונבולוציה וגם בעקיפין כפול פרמטרים נלמדים בשכבת ה- FC , מצליח כבר ליצור גרף למידה יורד, אך לא בהצלחה גבוהה.

בגרף הירוק ($conv1 = 12, conv2 = 32$), אנחנו כבר מתחילים להתייצב ולקבל תוצאה דומה לאדום בטסט למרות שאנחנו עם כפול פרמטרים ועם לוס נמוך יותר באימון. הגרף הסגול הראשון ($conv1 = 24, conv2 = 64$), שמתחיל להראות $overfit$ אבל מאוד קטן כשהלוס שלו בטסט קצת עולה מעל קודמיו בסוף. הכתום ($conv1 = 24, conv2 = 64$), מגיע כבר ללוס די גדול. עוקף אפילו את הכחול שכמעט ולא לומד.

והתכלת, שהוא כבר עם כמות נכבדת של פרמטרים, מצליח לשנן את הדאטא די ברצינות ולהגיע לתוצאות גרועות אפילו יותר בטסט לוס כשהאימון שלו ממש מתחת לכולם.

מסקנה כרגע שהפתרון האידיאלי מהשישה הוא הגרף הירוק ($conv1 = 12, conv2 = 32$), כי הוא היחיד שלא מגיע ל- $overfit$.

לצורך סגירת מעגל נחשב זריז את הפרמטרים שלו:

$conv1$: גודל כל קרנל הוא 5 על 5 על 3 שזה 75 ועוד באייס שזה 76. יש לנו 12 ערוצים כאלה ולכן: 912.

$conv2$: גודל כל קרנל הוא 5 על 5 על 12 שזה 300 ועוד 1 זה 301. יש לנו 32 ערוצים כאלה ולכן: 9632.

fc : יש לנו קלט בגודל 5 על 5 על 32 שכשנשטח אותו נקבל ווקטור באורך 800. הפלט שלנו הוא ווקטור באורך 10 ולכן צריך מטריצה באורך 800 על 10 אז 8000 ויש באייס של 10 ולכן 8010.

נסכום הכל ונקבל $11354 = 912 + 9632 + 8010$ פרמטרים נלמדים.

שאלה פרקטית שנייה Importance of Non – Linearity:

התבקשתי למחוק את הפעולות הלא לינאריות שברשת.

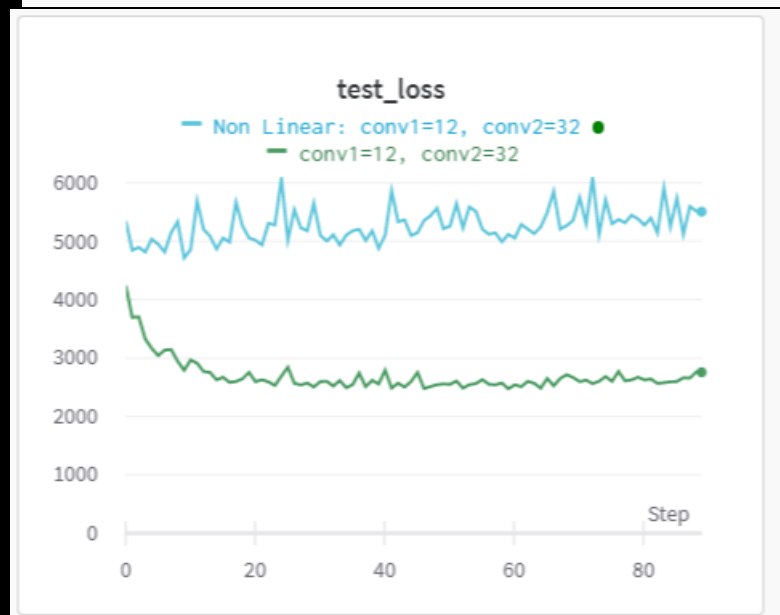
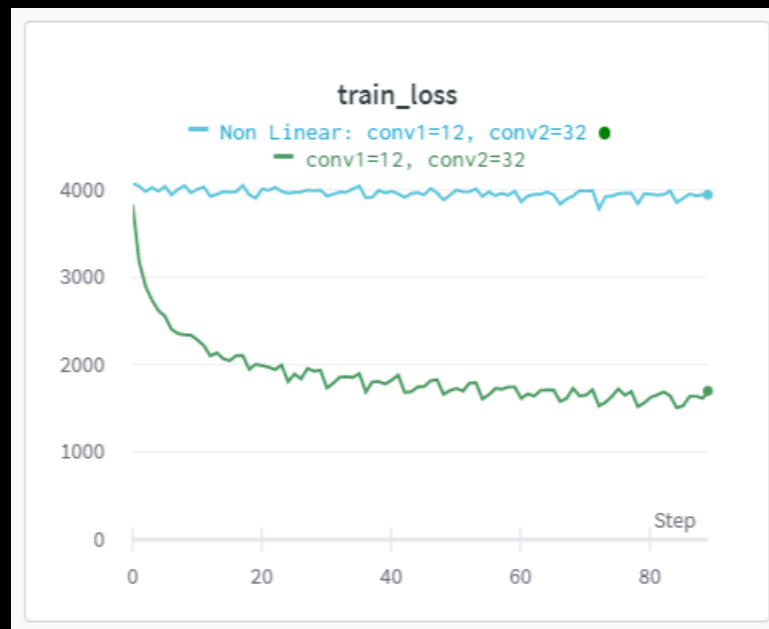
אז את האקטיבציות של $relu$ זה קל כי לא זה לא הורס את המימדים.

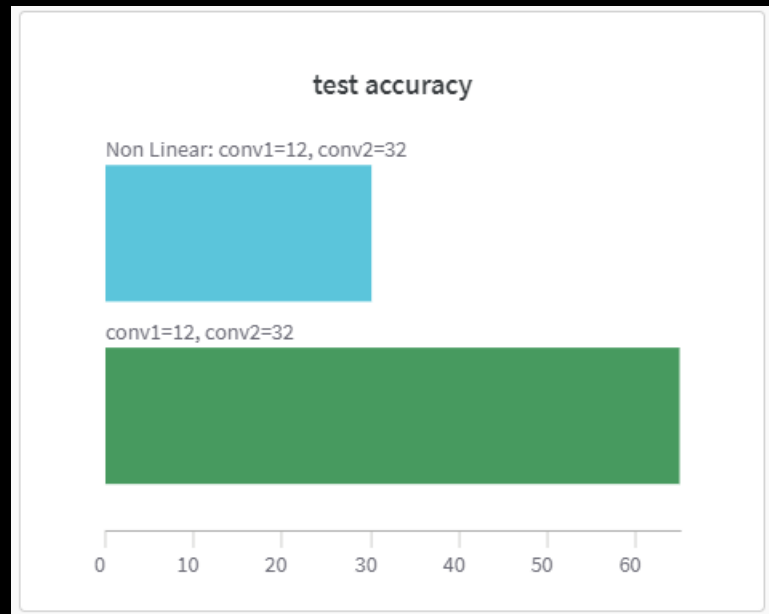
כשנוריד את $max pooling$ שזה גם לא לינארית אנחנו נקבל קלט לשכבת ה- fc שהיא בגודל 24 על 24 על 32 וזה 18,432 פרמטרים נלמדים רק לשכבת ה- fc .

זה מכיוון שלא צימצמנו את המימדים של הרשת.

לסיכום נבדוק עכשיו רשת עם 2 שכבות קונבולוציה בעלות קרנל 5 על 5 כשלראשונה פלט בעומק 12 ולשנייה פלט בעומק 32 והתוצאה תיכנס ל- fc פשוט שמחזיר וקטור באורך 10.

התוצאות של הרשת הזו הן (בתכלת התוצאות כשמורידים כל רכיב לא לינארי וירוק זו הרשת המקורית):



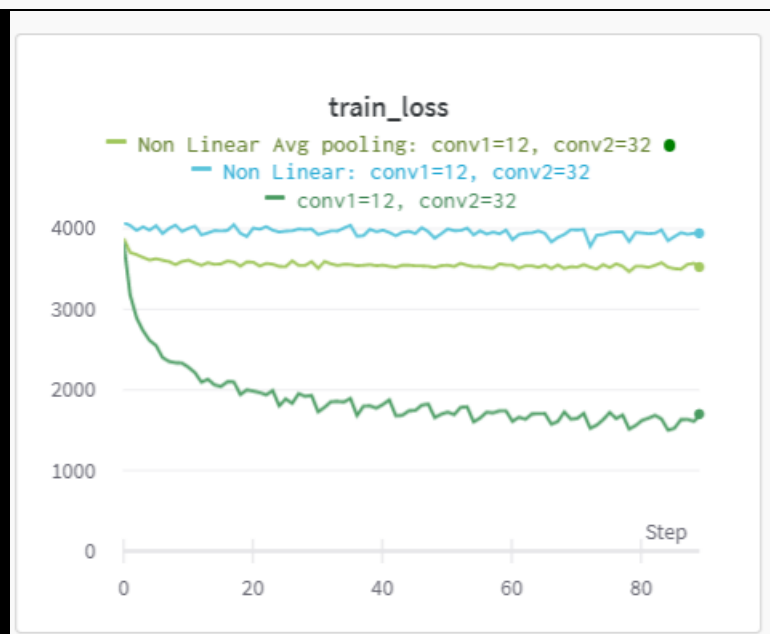
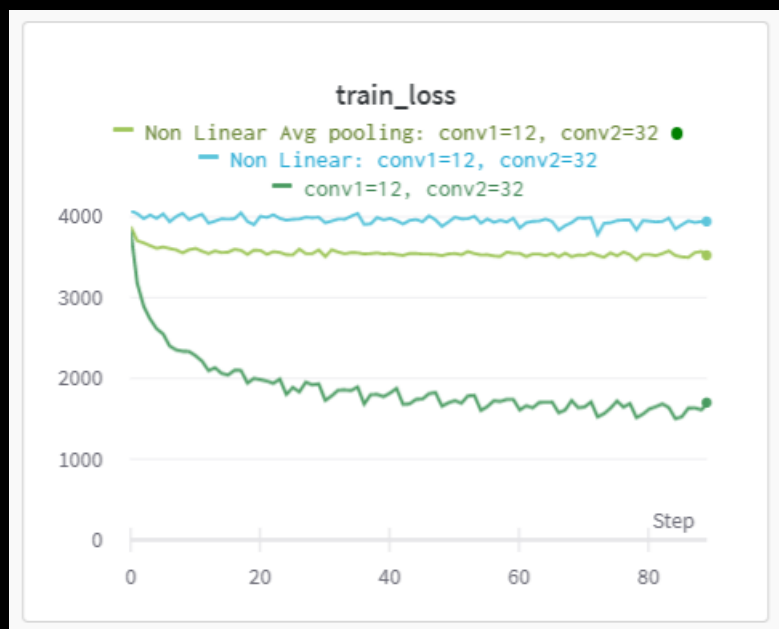


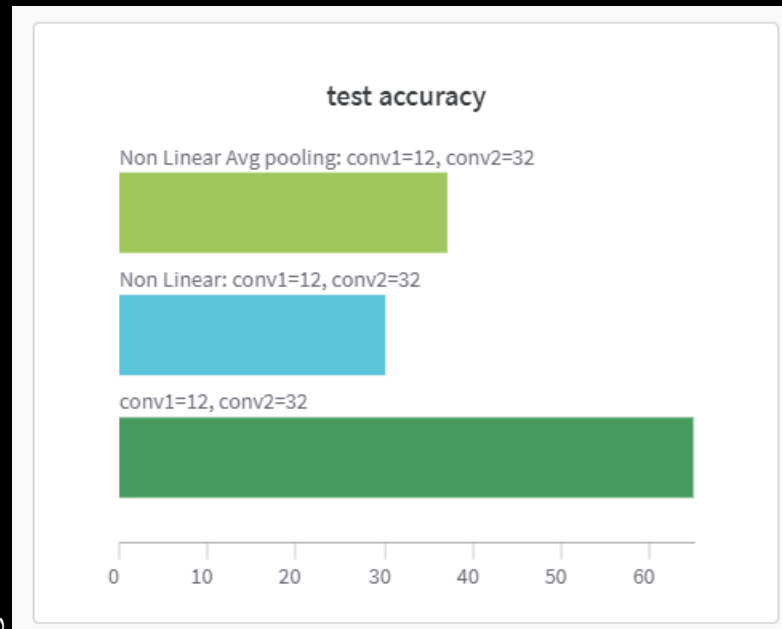
ניתן לראות שאפילו ה-*train loss* לא יורד וגם ה-*test loss* לא יציב ובטח שלא יורד והדיוק בהתאם, 30 אחוז, לעומת 65 במודל שברחנו מהסעיף הקודם.

בואו ננסה, לפני שנעביר לסעיף הבא לאפשר לרשת להיראות דומה יותר בכמות הפרמטרים למודל הירוק, כי יותר פרמטרים זה משימה אחרת שאולי לא מצליחה להתכנס פשוט.

במקרה זה נמיר את שכבות ה-*max pooling* לשכבות ה-*average pooling* וככה נישאר על אותה כמות פרמטרים בלי לינאריות.

אלו התוצאות: (הטורקיז זה הגרף החדש כתוספת על הגרפים הקודמים)





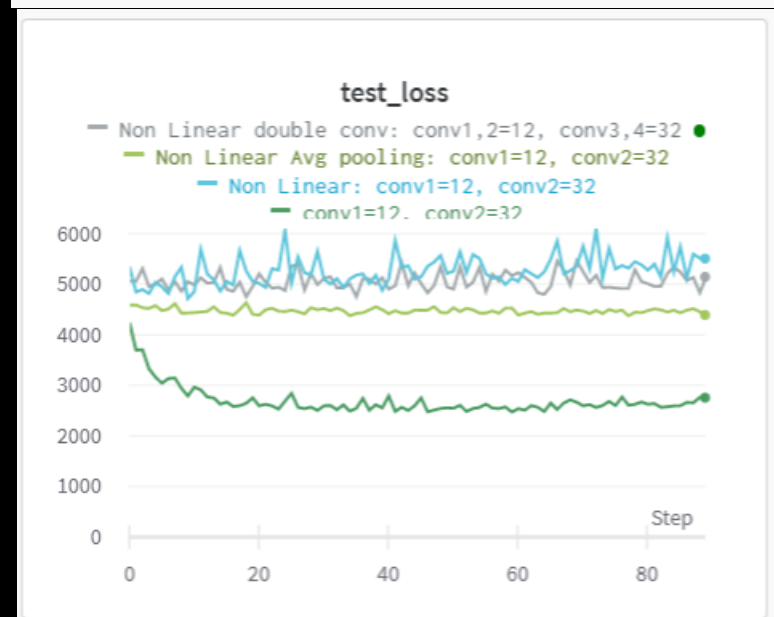
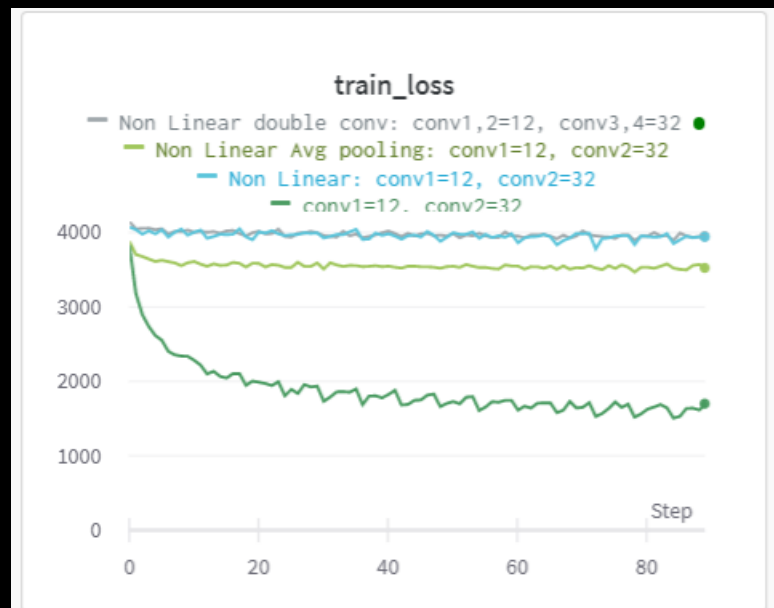
(באחוזים)

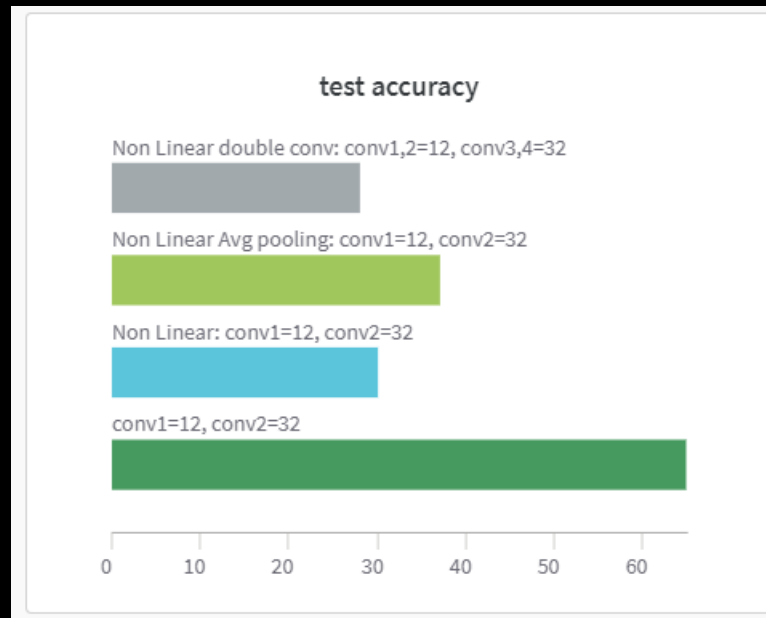
אפשר לראות שאפילו עכשיו, כשהשארנו את השינויים בגודל הקלט והפלט של כל שכבה באמצעות הוספת *avgPool* עדיין התוצאות משמעותיות שונות כשיש לינאריות.

לסיום נחזור לפתרון ללא *AvgPool* רק שהפעם יהיו לנו הרבה יותר שכבות קונבולוציה של 5 על 5.

אם עד עכשיו היו לנו 2 שכבות, הפעם נוסיף 2 שכבות ונקבל ארבע שכבות קונבולוציה. 2 ראשונות יהיו עם 12 ערוצים והאחרונות עם 32 ערוצים.

נרצה לקבל תוצאות דומות לגרף התכלת, כי כשאין לינאריות אז הכל זה כפל מטריצות אחד גדול. (גילוי נאות - שורה זו נכתבה לפני הרצת המודל ואני שמח שזה אכן מה שקרה) הנה התוצאות: (באפור)

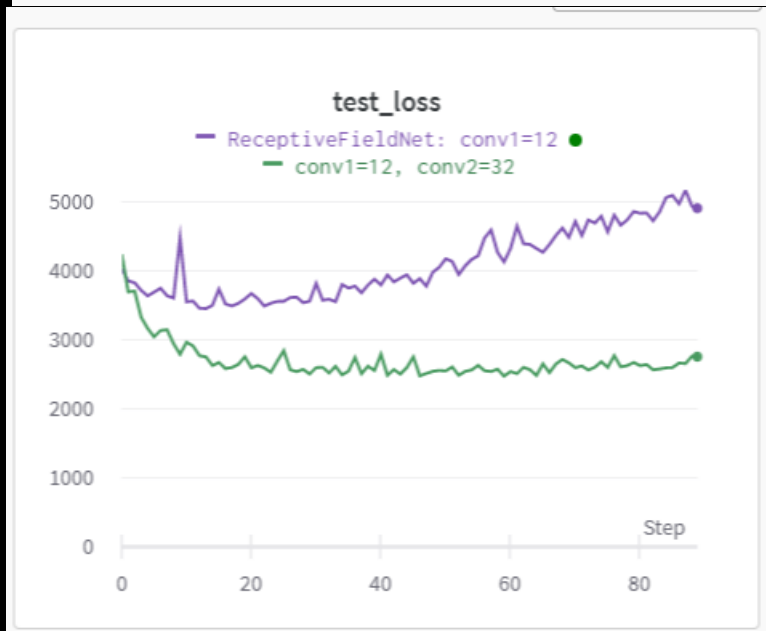
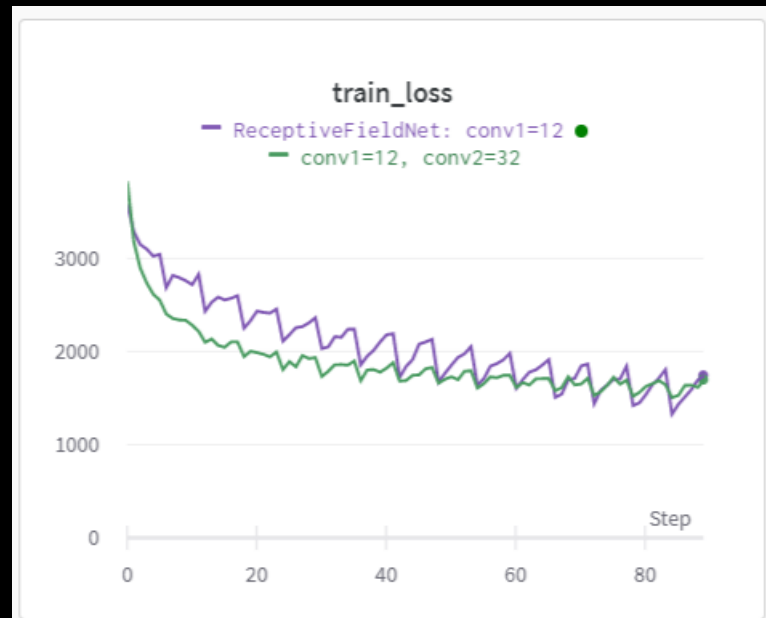


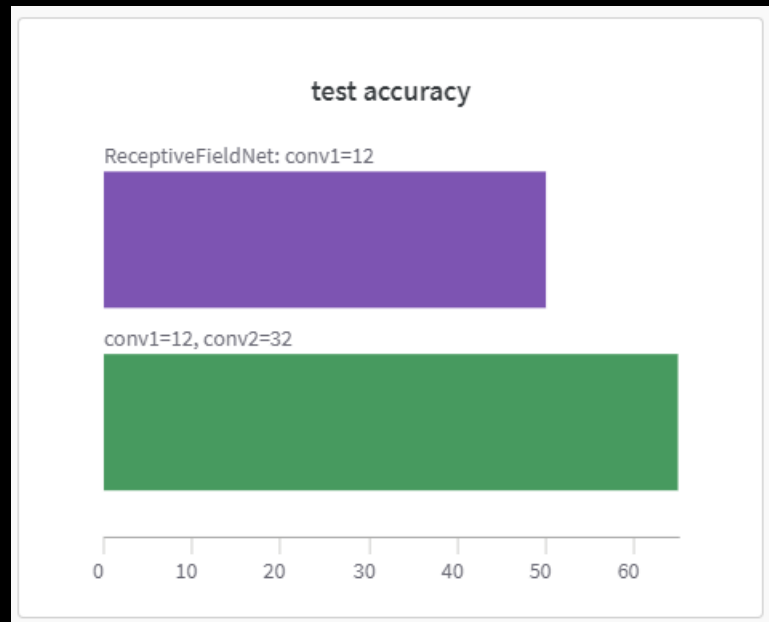


אפשר לראות בגרפים כמה מתואמים הגרפים האפורים והטורקיז, כשהבדל ביניהם היא מכפלה של כמות הפרמטרים הנלמדים בשכבות הקונבולוציה! כשהיה לנו לינאריות אפילו הזזה של פחות פרמטרים שינתה והפרידה בין הגרפים, במיוחד על סט האימון. כמו שכבר אמרתי, מאחר ויש לינאריות בכל הפרמטרים, אנחנו מקבלים תלות בין הפרמטרים ככה שהכל מסתכם למטריצה אחת באמצעות הכפלת מטריצות במקום למידה של מטריצות בלתי תלויות זו בזו.

שאלה פרקטית שלישית. *Cascaded Receptive Field*:

בחרתי באפשרות הראשונה ליישום, בעיקר כי היא נראתה לי כרשת שמציגה בקלות יותר את מטרת הסעיף - להראות את הכוח של *receptive field*. אז עכשיו יהיה לנו רשת עם שכבת קונבולוציה אחת עם קרנל 5 על 5 ערוצים שאת הפלט שלה נכניס לשכבת ה-*fc* שלנו שתחזיר פלט באורך 10. הנה התוצאות בסגול ביחס לתוצאות המקוריות בירוק:





נבחין בשגיאה הגבוהה ב- $test\ loss$ שאמור להעיד על $over\ fit$ כי השגיאה של $train$ כן ירדה, כנראה כי לא היה הרבה פרמטרים ללמוד. מצד שני הדיוק הוא 50 אחוז, יחסית גרוע, יותר טוב מהטלת מטבע כי יש לנו 10 מחלקות אבל עדיין ראינו שהרבה מהמודלים הפשוטים שלנו הצליחו לעקוף את זה. נראה שהבעיה היא לא במספר הפרמטרים ללמידה: בשכבת $conv1$: היו לנו קרנלים בגודל 5 על 5 על 3 שזה 75 עם באייס 76 פרמטרים לכל קרנל כשיש לנו 12 קרנלים אז $912 = 76 * 12$ ובשכבה fc : מקבלים קלט בגודל $9408 = 28 * 28 * 12$ ופלט בגודל 10 ולכן עם הבאייס: $9409 * 10 = 94090$. שזה כמות רצינית של פרמטרים ללמידה. ולכן הבעיה כאן היא שלא הצלחנו בכלל לתפוס כמו שצריך את התמונה ולכן הלמידה לא הועילה.

שאלות תיאורטיות:

(1) נסביר את השאלה: אנחנו רוצים להוכיח שאם מתקיים:

$$L[x(i+k)]_{(j)} = L[x(i)]_{(j+k)}$$

כש L הוא אופרטור לינארי, x הוא וקטור קלט, $x(i)$ זה הקוארדינטה ה- i בוקטור הקלט, k הוא סקלר שיסמן לנו אופסט מסוים על וקטור הקלט. $L[x(i)]$ מחזיר לנו וקטור פלט מהפעולה L על סקלר כלשהו. L היא הכפלה במטריצה ולכן לכל קוארדינטה בוקטור הקלט אנחנו מקבלים וקטור כפלט. הפעולה $L[x(i)]_{(j)}$ מסמנת את הקוארדינטה ה- j בוקטור הפלט. עלינו להוכיח שמדובר פה בפעולת קונבולוציה. רמז: נפרק את $x(i)$ לסכום ממושקל של פונקציות δ ואז להשתמש בלינאריות של L . שאלה מכווינה: איזה סיגנל כקלט ייתן לנו פעולה דומה לקונבולוציה?

פתרון: נראה שמדובר באופרטור קונבולוציה.
נבחין כי:

$$x(i) = \sum_{k=1}^n x(k) \cdot \delta(i - k)$$

מלינאריות L נקבל כי:

$$L[x(i)]_{(j)} = L \left[\sum_{k=1}^n x(k) \cdot \delta(i - k) \right]_{(j)} = \sum_{k=1}^n [x(k) \cdot L[\delta(i - k)]_{(j)}]$$

$$= \sum_{k=1}^n [x(k) \cdot L[\delta(i)]_{(j-k)}] = \sum_{k=1}^n [x(k) \cdot l(j - k)] = (l * x)_j$$

(2) אין משמעות לסדר, מכיוון שאנחנו הולכים להכפיל את הוקטור תוצאה במטריצה שאותה נלמד, אז כל ערך יחיד i בוקטור הולך להיות מוכפל בשורה ייעודית i במטריצה ולכן אנחנו פשוט נלמד את השורה הזאת להתאים לערכים שמגיעים לקואורדינטה i בקלט. ואם היינו מחליפים את i ו- j , מכיוון שאין תלות בין המשקולות ברשת אז היינו מקבלים שהשורות המתאימות היו מתחלפות במטריצת המשקולות פשוט. (הערה: הבנתי שיש מקרים בהם עבודה עם *batch* מייצרת איזשהי תלות בסדר ברמה הפרקטית ולא התיאורטית, אבל אני לא מכיר את זה מספיק ולכן אני משאיר את התשובה שלי כמו שהיא).

(3 א) לא, *Relu* אינה *LTI* מכיוון שהיא אינה לינארית. אפשר לקחת דוגמא קטנה של הזזה של ערך יחיד 1 בקלט במינוס 2 ולראות שהערך יגיע ל-0 ואם נזיז אחרי ההפעלה נגיע ל-1.

(ב) לא, *strided pooling* אמנם פעולה לינארית ואפשר להציג אותה כמטריצה מתאימה עם עמודות 1 ו-0. אך אם נבצע *translation* על הקלט הפלט יכול לקבל שינוי חד בערכים לפי המיקום של הערכים במיקומים המתאימים.

(ג) לא, מכיוון שחיבור בקבוע אינו לינארי. אפשר לקחת וקטור פשוט לבצע הזזה עליו ולראות שהסכום החדש ייראה שונה לחלוטין.

(ד) לא, אמנם יש לנו כאן הכפלה במטריצה אבל בהנחה ונזיז את וקטור קלט אז הערכים יוכפלו בעמודות שונות שיכולות להיות עם ערכים שונים זה מזה כך שהפלט ייצא שונה.