

## Project 2 - CNN

מגישה: גל חלילי

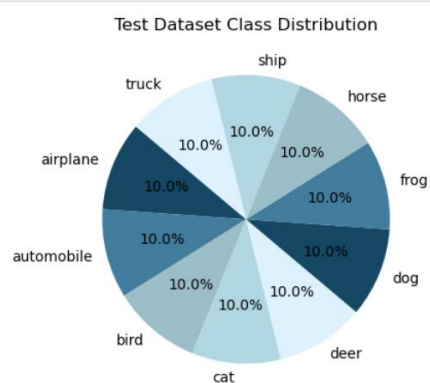
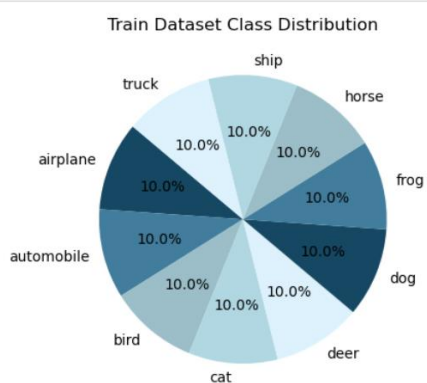
Data exploration

Dataset labels - classes: ['airplane', 'automobile', 'bird', 'cat', 'deer', 'dog', 'frog', 'horse', 'ship', 'truck']

First 5 samples from Train Dataset



First 5 samples from Test Dataset



\* ניתן לראות שכל המחלקות מתפלגות באופן אחיד גם בtrain וגם בtest – בשני הסטים בכל מחלקה יש 10% מדוגמאות הסט.

בניית הרשת:מודל 1 –

התחלתי בניסוי הראשון עם שתי שכבות קונבולוציה ואחריהן 4 שכבות ליניאריות FC.

שכבת הקונבולוציה הראשונה לוקחת את התמונות ומפעילה עליהן 6 פילטרים בגודל 5X5. השכבה השנייה לוקחת את הפלט של השכבה הקודמת (שזו בעצם activation map עם 6 ערוצים) ועליו מפעילה 12 פילטרים בגודל 5X5.

השכבה השלישית היא בעצם שכבת ה-FC הראשונה – יש בה 120 ניוונים, אחריה שכבה עם 60 ניוונים, לאחר מכן שכבה עם 40 ניוונים ולבסוף שכבת הפלט עם 10 ניוונים בהתאמה לכך שיש 10 מחלקות בבעיה זו.

בין כל שכבה ושכבה מופעלת פונקציית אקטיבציה RelU (בחרתי בה כי בהרצאות אמרנו שלרוב היא נותנת תוצאות טובות). בנוסף, פונקציית הלוס היא cross-entropy ושיטת האופטימיזציה היא SGD עם  $lr=0.001$ .

אחרי כל שכבת קונבולוציה הפעלתי שכבת Max-Pooling:

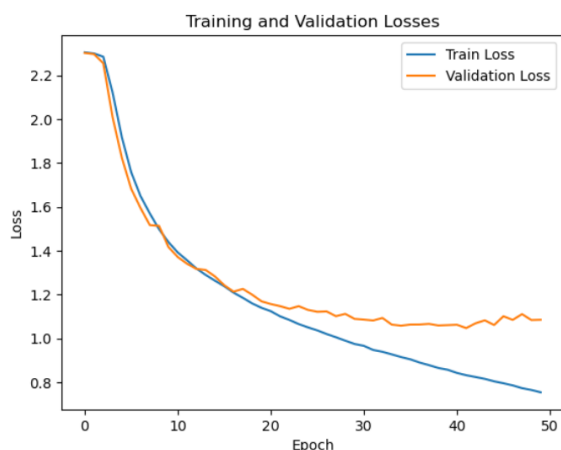
זו שכבה לא למידה שמטרתה להקטין את המימד של כל activation map בנפרד. זה סוג של פילטר שמחושב כמו פילטר רגיל אבל לוקח את הערך המקסימלי מכל חלון ולא סכום (כלומר עושה את המכפלות בין כל ערך בפילטר לבין הערך המתאים בחלון התמונה ובמרום לכלום את המכפלות, לוקח את המכפלה המקסימלית) – כך העצם מקטינים את מימד התמונה כי במקום לשמור את כל המשקולות, שומרים רק את המשקולות של הפיצאר הכי משמעותי מכל חלון.

בשכבות אלו לקחתי פילטר בגודל 2X2 עם stride 2.

\*בשביל להכניס את הפלט של שכבת max pooling האחרונה לתוך שכבת ה FC הראשונה צריך "לשטח" את הטנזור לוקטור (חד מימדי).

\*אחרי שכבת הפלט לא מופעלת RelU כי שם מופעלת softmax בשביל לקבל את ההסתברויות של להיות שייך לכל מחלקה.

גרף שמציג את loss של סט האימון וסט הוולידציה:



הגעתי עם מודל זה ל  $accuracy = 0.6373$  בוולידציה שזה לא

דיוק מספק. בנוסף ניתן לראות פה overfitting.

## מודל 2 –

אחרי שראיתי שנוצר במודל הראשון overfit ניסיתי להשתמש ב dropout ו-early stopping בכדי למנוע זאת.

### -Dropout

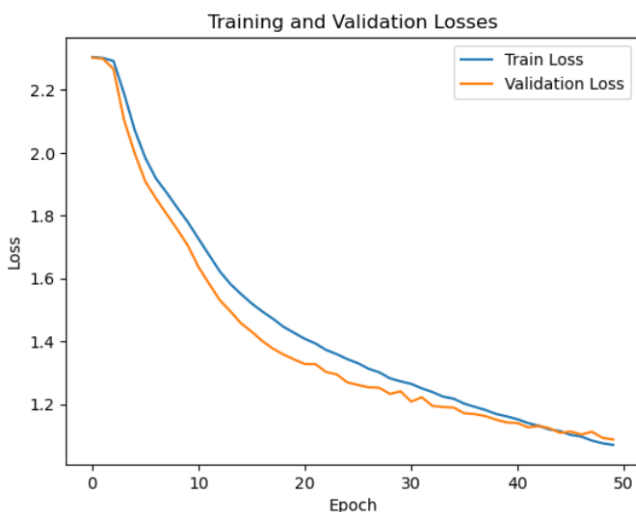
זו שכבה שמטרתה לאפס חלק מהנוירונים משמע לא להעביר דרכה את הערכים. מטרת השכבה ליצור הכללה ולמנוע overfitting. בכל איטרציה אנחנו משמיטים חלק מהנוירונים, נאפס כל נוירון בסיכוי  $p$  (  $p$  גדול יותר גורר יותר נוירונים מאופסים) כך שבעצם חסר לו מידע אחר כל פעם והוא עדיין מנסה להצליח (למזער את הלוס), מה שגורם לו "לא להיות תלוי" בקומבינציה הספציפית של הפיצ'רים שראה בסט האימון – כלומר לומד להכליל ולא לשנן את הקשרים הספציפיים שראה.

### -Early stopping

גם מספר epochs יכול להשפיע על overfit. נסתכל על ה-validation set loss (בקיצור vsl) ונחליט לעצור את האימון כך שהוא לא יגדל. בעצם השיטה אומרת שכל פעם שמזהים עלייה ב vsl נשמור את ה epoch בו זה קרה ונתחיל לספור  $n$  (לבחירתנו) epochs שבהם נבדוק אם זו אכן מגמת עלייה אמיתית או שזו הייתה עלייה לצורך ירידה. אם זו אכן מכמת עלייה ולא זיהינו מאז ירידה ב vsl אז נחליט לעצור את הריצה ב epoch ששמרנו. אבל אם זיהינו מאז ירידה אז נמשיך בתהליך.

שיטה זו עוזרת במניעת overfitting בגלל שעליה ב loss של ה-validation מסמנת על overfit לכן לפי עלייה זו נחליט מתי לעצור בשביל לא להיכנס ל overfitting. בגלל ששיטה זו עוזרת למנוע overfitting היא עוזרת למודל להכליל טוב יותר במקום לשנן כי עוצרת אותו לפני שהוא מתחיל ללמוד קשרים ותלויות שספציפיים לדוגמאות שהוא רואה בסט האימון ואם ילמד וישנן אותן הוא יהיה טוב מאוד על ה train ופחות טוב על ה validation ולכן vsl יעלה.

השארתי את כל המודל כפי שהיה ורק הוספתי אחרי כל שכבת FC (חוץ משכבת הפלט כמובן) שכבת dropout שבה ה- dropout rate  $p=0.2$ . בנוסף, קבעתי שמספר האפוקים המקסימלי יהיה 100 וש-  $n$  (מספר האפוקים שבהם מחכים להשתפרות בלוס של הוולידציה) יהיה 10.



במודל זה הצלחתי להגיע ל – validation acc = 0.618  
הצלחתי למנוע את ה overfitting בצורה משמעותית  
אבל עדיין המודל לא לומד מספיק והביצועים נמוכים.

יכול להיות שהרשת שהגדרתי לא מספיק מורכבת בשביל הבעיה הזו וצריך יותר שכבות ויותר פילטרים בשביל לשפר את הלמידה.

### מודל 3-

במודל זה שמתי 3 שכבות קונבולוציה ו-3 שכבות FC.

שכבות הקונבולוציה:

בכל אחת גודל הפילטר הוא  $3 \times 3$  עם padding 1.  
בשכבה הראשונה מופעלים 32 פילטרים, בשכבה השנייה 64 ובשלישית 128.

אחרי כל שכבת קונבולוציה שמתי שכבת BatchNorm2d:

Batch Norm is a normalization technique done between the layers of a Neural Network instead of in the raw data. It is done along mini-batches instead of the full data set. It serves to speed up training and use higher learning rates, making learning easier.

לא למדנו על שיטה זו בהרצאות אך ממה שהבנתי מההסברים שקראתי (וצירפתי אחד מהם למעלה) זה ששכבה זו אמורה לסייע בייצוב והאצת תהליך האימון על ידי נרמול הקלט לכל שכבה.

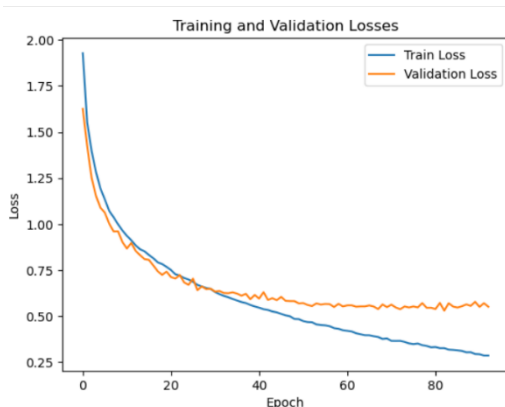
שכבות ה-FC:

בשכבה הראשונה יש 256 ניוונים, בשנייה 64 ובשכבה האחרונה-שכבת הפלט 10 ניוונים. בהתאם לכך שיש 10 קלאסים.

שכבת ה-FC הראשונה מקבלת כקלט את הפלט משכבת הקונבולוציה אחרי שנהפך לוקטור (חד מימדי).

גם פה השתמשתי ב max pooling בגודל  $2 \times 2$  ו-stride 2 אחרי כל שכבת קונבולוציה (אחרי BatchNorm2d). ושוב פונקציית האקטיבציה הייתה ReLU עם  $lr=0.001$ .

\*השארתי את שיטות מניעת ה-overfitting שהשתמשתי בהן במודל הקודם ( $p=0.2$ ,  $n=10$ ).



במודל זה validation accuracy = 0.8179

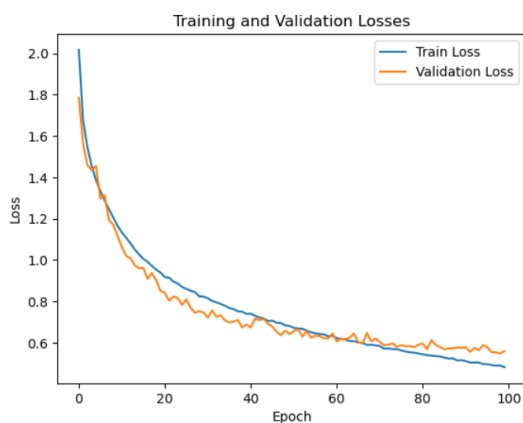
אכן יש שיפור משמעותי בביצועים (כמעט 20%)

אבל שוב נוצר overfitting.

יכול להיות שעכשיו כשהרשת עמוקה יותר הלמידה נהייתה גם כן עמוקה וספציפית יותר וקו שערור במודל 2 למנוע את ה-overfitting כעת לא ממש עוזרים. אנסה להגדיל את  $p$  ל-0.3.

#### מודל 4 -

מודל זה זהה למודל הקודם – רק הגדלתי את  $p$  ל-0.3.



במודל זה validation accuracy = 0.8123

הגדלת  $p$  אכן עזרה במניעת overfitting אך עדיין יש שוני בין הלוס של האימון לוולידציה.

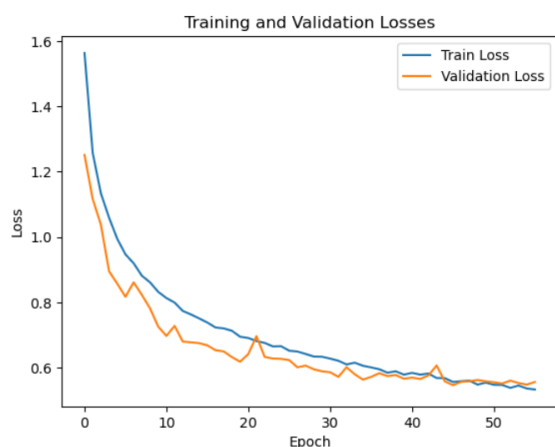
מה שניכר מגרף זה לדעתי זה שהגרף של הוולידציה מאוד רועש – אולי כדאי לשנות את שיטת האופטימיזציה.

מודל 5 -

שוב השתמשתי באותה רשת רק שהפעם האופטימיזר הוא Adam.

זו שיטת אופטימיזציה שלא למדנו עליה בקורס אך ממה שקראתי והבנתי: Adaptive Moment Estimation, הוא אלגוריתם קצב למידה אדפטיבי שנועד לשפר את מהירויות האימון ברשתות עצביות עמוקות ולהגיע להתכנסות במהירות. הוא מתאים אישית את קצב הלמידה של כל פרמטר על סמך היסטוריית השיפועים שלו, והתאמה זו עוזרת לרשת העצבית ללמוד ביעילות.

באלגוריתמים שאנחנו למדנו, קצב הלמידה  $\alpha$  קבוע (hyperparameter), כלומר עלינו להתחיל בקצב למידה כשלהו ולשנות ידנית את האלפא לפי שלבים או לפי לוח זמנים למידה כלשהו. קצב למידה נמוך יותר בתחילתו יוביל להתכנסות איטית מאוד, בעוד ששיעור גבוה מאוד בהתחלה עלול להחמיץ את המינימום. Adam פותר בעיה זו על ידי התאמת קצב הלמידה  $\alpha$  לכל פרמטר  $\theta$ , מה שמאפשר התכנסות מהירה יותר בהשוואה לגרדיאנט דיסנט סטנדרטי עם קצב למידה קבוע.



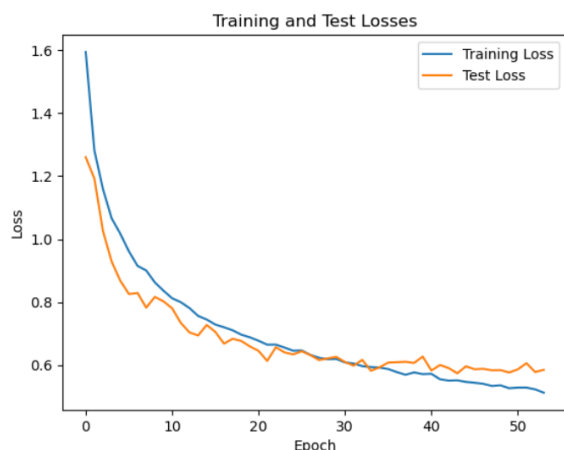
במודל זה validation accuracy = 0.8131

הגרף אכן פחות רועש מהמודל הקודם וגם ה-overfitting שהיה ירד.

בנוסף, ניתן לראות שאכן התכנסנו הרבה יותר מהר במודל הקודם- במודל 4 היו 100 אפוקים ואילו במודל הנוכחי קצת יותר מ-50.

כלומר המודל היה צריך קרוב לחצי מכמות האפוקים במודל 4!

בחרתי במודל האחרון וכעת אני אאמן את המודל על כל סט האימון (ללא וולידציה) ואז אשערך על test set-ה.



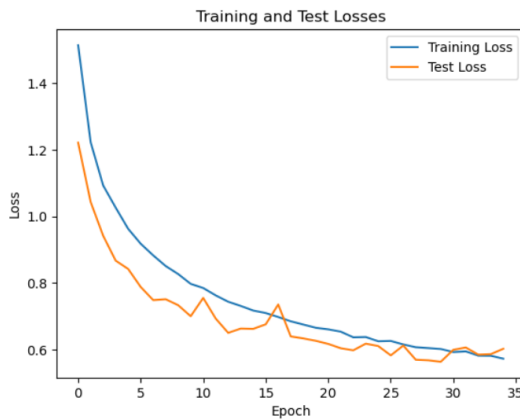
קיבלתי Test Accuracy = 0.8068

יש עדיין קצת overfitting בסוף.

אולי המודל למד קשר כלשהו שספציפי לדוגמאות מהאימון

והוולידציה ועכשיו בtest יש קצת חוסר התאמה ולכן יש overfitting.

אני אנסה להריץ את אותו מודל ורק להקטין את n ל-5 כי אולי אני נותנת יותר מידי זמן השתפרות.



אכן ההקטנה של n עזרה להוריד את overfitting.

והצלחתי להגיע ל-0.8065 Test Accuracy

### הסבר מימדים של הרשת שנבחרה:

$32 \times 32 \times 3 \rightarrow \text{RelU}(\text{BatchNorm2d}(\text{Conv}(\text{in: } 3, \text{ out: } 32, \text{ kernel: } 3, \text{ paddind: } 1) \rightarrow (32+2*1-3)/1 + 1 \rightarrow 32 \times 32 \times 32$

$32 \times 32 \times 32 \rightarrow \text{Max-Pooling}(\text{kernel: } 2, \text{ stride: } 2) \rightarrow (32-2)/2 + 1 \rightarrow 16 \times 16 \times 32$

Dropout()

$16 \times 16 \times 32 \rightarrow \text{RelU}(\text{BatchNorm2d}(\text{Conv}(\text{in: } 32, \text{ out: } 64, \text{ kernel: } 3, \text{ paddind: } 1) \rightarrow (16+2*1-3)/1 + 1 \rightarrow 16 \times 16 \times 64$

$16 \times 16 \times 64 \rightarrow \text{Max-Pooling}(\text{kernel: } 2, \text{ stride: } 2) \rightarrow (16-2)/2 + 1 \rightarrow 8 \times 8 \times 64$

Dropout()

$8 \times 8 \times 64 \rightarrow \text{RelU}(\text{BatchNorm2d}(\text{Conv}(\text{in: } 64, \text{ out: } 128, \text{ kernel: } 3, \text{ paddind: } 1) \rightarrow (8+2*1-3)/1 + 1 \rightarrow 8 \times 8 \times 128$

$8 \times 8 \times 128 \rightarrow \text{Max-Pooling}(\text{kernel: } 2, \text{ stride: } 2) \rightarrow (8-2)/2 + 1 \rightarrow 4 \times 4 \times 128$

Dropout()

Flatten  $4 \times 4 \times 128$  to  $4*4*128$

Dropout(ReLU(Linear(in:  $4*4*128$ , out: 256)

Dropout(ReLU(Linear(in: 256, out: 128)

Linear(in: 128, out: 10)