

Assignment 4 – implementing an SNLI paper

במשימה התבקשנו לממש מודל שמישיתו היא להכריע עבור 2 משפטים - premise, hypothesis, האם משפט אחד נובע מהשני (entailment), האם המשפטים סותרים (contradiction), או האם אין קשר ביניהם (neutral).

עבור משימה זו נכתבו מאמרים שונים שמתארים מודלים שונים שהשתתפו בתחרות של סטנפורד על הדאטה שהם ייצרו – SNLI corpus.

בתחילה בחרנו מאמר בשם "600D (300+300) Deep Gated Attn. BiLSTM encoders" שנוקט בשיטה של להפוך כל משפט לווקטור בגודל קבוע ומשם לעבור לסיווג. מכיוון ששיטה זו משתמשת ברשתות עמוקות עם הרבה פרמטרים לאימון, הבנו שתהליך האימון יהיה ארוך מדיי וכבד מדיי ביחס לזמן ולמשאבים שיש לנו. לכן החלטנו לבחור במאמר אחר: "A Decomposable Attention Model for Natural Language Inference", vanilla approach.

הרעיון המרכזי במאמר הוא להתייחס למבנה המשפטים עם תשומת לב לנושא ולפועל, ולמעשה לבצע soft alignment בין hypothesis לpremise באמצעות attention.

תיאור המודל:

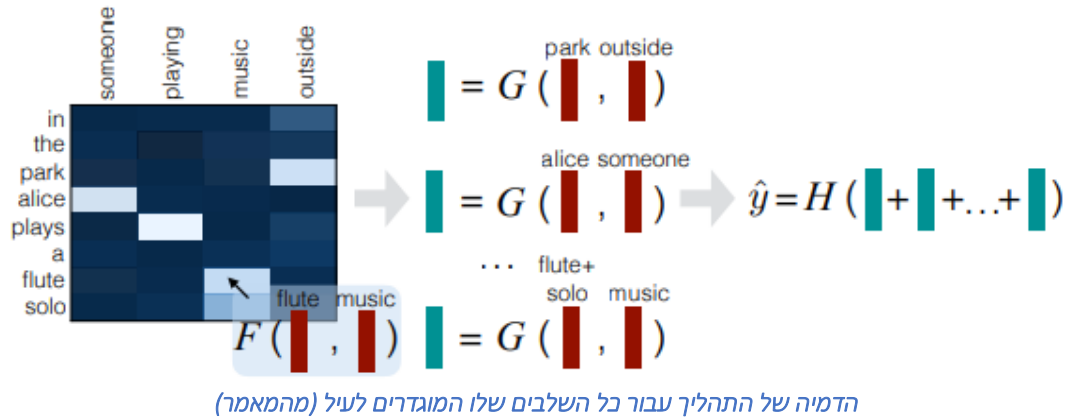
כדי לממש את המאמר בצורה טובה ביותר פעלנו כך:

1. טעינת הדאטה - torchtext:

- כדי לטעון את הדאטה השתמשנו בספריית torchtext שטוענת אוטומטית את כל הדאטה מתוך הקבצים של train, dev, test אל טנזורים בפייתון שאפשר להשתמש בהם בנוחות. באמצעות הדגל Pre_processing הוספנו לכל משפט NULL בתחילתו כפי שהתבקש במאמר.
- על הדאטה שטענו בנינו bucket_iterator שממין את הדוגמאות לפי אורכי המשפטים מהקצר אל הארוך כך שנדרש פחות padding בכל batch ובעקבות כך זמן העיבוד מתקצר משמעותית.
- על בסיס הקבצים שנטענו בנינו את אוצר המילים עבור טבלת embedding וכל מילה קיבלה אינדקס. (הסבר מפורט על שיטת embedding בהמשך)
- עבור כל דוגמה בdata_set נבצע את השלבים הבאים:

- שלב ה-ENCODE:** כל מילה בכל משפט (premise, hypothesis) קיבלה את הייצוג שלה לפי טבלת embedding. את הווקטורים מטבלת embedding הגדרנו להיות במידה 300 והעברנו אותם תחת שכבה לינארית כדי להקטין את גודלם ל-200.
- שלב ה-ATTEND:** מחשבים עבור כל מילה במשפט הראשון את המשקלים אל מול כל המילים במשפט השני (באמצעות מעבר ברשת נירונים בעלת 2 שכבות, פונקציית אקטיבציה relu). לאחר מכן אנו מבצעים כפל בין המשקלים למשפטים המקוריים ויוצרים את שני הווקטורים beta1 alpha ושולחים לשלב הבא. למעשה מה שהתבצע הוא soft_align בין שני המשפטים הנתונים.
- שלב ה-COMPARE:** בהינתן וקטורי beta1 alpha נשרשר אותם עם המשפטים המקוריים בהתאמה ואת התוצאה מעבירים בעוד שכבת נירונים.
- שלב ה-AGGREGATE:** בהינתן וקטורי v1, v2 מהשלב הקודם נסכום כל אחד מהווקטורים ונשרשר אותם יחד, את השירשור נעביר ברשת נירונים.
- את התוצאה שקיבלנו מפונקציית aggregate נעביר בשכבה לינארית נוספת שהoutput שלה הוא וקטור בגודל מספר הlabels שלנו.

- הכניסה שקיבלה את הציון הגבוה ביותר נבחרה להיות הפרדיקציה.



הסיבה שבחרנו במאמר: כמות הפרמטרים קטנה יחסית והמשאבים שלנו מאפשרים את האימון הנדרש. בנוסף, ההתייחסות למבנה המשפט ולא רק למשמעות המילים נתנה פרספקטיבה חדשה על איך להתמודד עם המשימה, ועניין אותנו לשחזר את ההצלחה של מאמר זה.

הסבר על שכבת הEMBEDDING:

לאחר טעינת הדאטה יצרנו את אוצר המילים עבור טבלת הembedding באמצעות אותה ספרייה (torchtext). אל אוצר המילים הוספנו וקטורי Pre_embedding של glove. עבור כל מילה באוצר המילים שנבנה ניתן אינדקס בהתאם לטבלה שנטענה מראש. Padding קיבל את אינדקס 1, ומילים לא ידועות קיבלו אינדקס 0. בתחילה טענו את הווקטורים מקובץ glove.6B.300D. קובץ זה מתייחס למילים בlower_case ולכן כשטענו את המשפטים הפכנו את כולם לlower_case כדי להתאים לאוצר המילים. לאחר בדיקות באינטרנט גילינו שהקובץ glove.840B.300D יותר רחב וגם נותן התייחסות שונה לlower וupper מה שמאפשר למידה מדויקת יותר למשפטים שנצפים בדוגמאות. שינוי זה אכן שיפר משמעותית את אחוזי הדיוק על dev_set. את הווקטורים שנטענו נרמלנו שוב (max_norm=1, נרמול l2) כפי שנתבקשנו במאמר. בנוסף, טבלת הembedding לא נלמדה במהלך האימון והווקטורים נשארו כמו שהם.

:OOV

לפי הוראות המאמר, עבור כל מילה שלא נמצאת בטבלת word_embedding (וקיבלה אינדקס 0) הוגרל וקטור אקראי מבין 100 וקטורים מאותחלים נורמלית בגודל 300 ששורשרו לטבלת הword_embedding המצויה.

נקודות בתהליך שיחזור ביצועי המאמר:

1. Optimizer – על פי המאמר בחרנו את AdaGrad עם LR=0.05, initial_accumulator_value=0.1. אחד מניסיונות השיפור שלנו כלל להגדיר את Adam הפופולרי optimizer עם אותו קצב למידה, אך לא קיבלנו שיפור ונשארו עם מה שהוצע במאמר.
2. hidden_layers בכל אחת מהרשתות הן בגודל 200 חוץ מהשכבה האחרונה, שכבת הoutput.
3. לכל אחת מהרשתות ביצענו רגולריזציה של dropout בהסתברות של 0.2, dropout בוצע לפני כל הפעלת פונקציית האקטיבציה-reLU, אך לא עבור השכבה הליניארית האחרונה.

4. את כל הפרמטרים של השכבות הלינאריות נתבקשנו לאתחל באיתחול גאוסייני ($\text{std}=0.01, \text{mean}=0$), כשעשינו זאת ידנית האחוזים ירדו משמעותית וקצב הלמידה היה קטן. דווקא האתחול הדיפולטי של Pytorch השיג תוצאות טובות יותר ולכן נשארנו איתו (`init.uniform_`).
5. Batch- במאמר דווח שהשתמשו בbatch בגודל 4, מכיוון שתהליך האימון מאוד איטי עם גודל כזה, השתמשנו בbatch בגודל 32 (שגם הוא הומלץ במאמר).
6. בטעינת הדאטה בחרנו לדרוס את שיטת `tokenized` הדיפולטיבית עבור המשפטים ולהשתמש דווקא בספריית `spacy` שמפרקת את המשפטים באמצעות אלגוריתם מדויק יותר עבור כללי הדקדוק בשפה האנגלית. שינוי זה הביא לשיפור משמעותי של לפחות 3% על `dev_set`.
7. `Weight_decay` - עבור `optimizer` שבחרנו הוספנו רגולרזציה מסוג L2 לחישוב `loss` עם מקדם 0.000001 דבר ששיפר במעט את הדיוק.
8. `Epochs` - אימנו את המודל על כל קבוצת האימון במשך 300 אפוקים כאשר בכל אפוק חישבנו את הדיוק על כל `dev_set`. הבחנו כי בערך באפוק 170 אנחנו מתכנסים לאותם אחוזים ולכן בהרצה הסופית אימנו את הפרמטרים במשך 170 אפוקים שלאחריהם עברנו על `test` וחישבנו את הדיוק.

מידת ההצלחה:

במאמר דווחו 86.3 אחוזי דיוק על `test`, 89.5 על `train`. המודל שלנו הצליח לשחזר את ההצלחה עבור `test` ואף להתעלות- הגענו ל-86.5 אחוזים. לעומת זאת על `train` הגענו רק ל-85.3.

ייתכן כי הוספת הרגולרזציה ללוס איפשרה יותר טעויות מה שהוביל להכללה טובה יותר על `test`. אך בתמורה הדיוק על `train` ירד.

learning curve graphs:

