

סקירה זו היא חלק מפינה קבועה בה אני סוקר מאמרים חשובים בתחום ה-ML/DL, וכותב גרסה פשוטה וברורה יותר שלהם בעברית. במידה ותרצו לקרוא את המאמרים הנוספים שסיכמתי, אתם מוזמנים לבדוק את העמוד שמרכז אותם תחת השם deepnightlearners.

לילה טוב חברים, היום אנחנו שוב בפינתנו deepnightlearners עם סקירה של מאמר בתחום הלמידה העמוקה. היום בחרתי לסקירה את המאמר שנקרא:

Teaching with Commentaries

של ג'ף הינטון האגדי ושותפיו.

פינת הסוקר:

המלצת קריאה ממייד: מומלץ לאוהבי מטה-למידה ובעלי רקע בחדו"א 2 מתקדם.

בהירות כתיבה: בינונית.

רמת היכרות עם כלים מתמטיים וטכניקות של ML/DL הנדרשים להבנת מאמר: רקע טוב בתחום מטה-למידה, חדו"א ברמה גבוהה.

יישומים פרקטיים אפשריים: ניתן להשתמש בגישה זו למשל לזיהוי דוגמאות המשפיעות ביותר על הביצועים או איתור פאטצ'ים בתמונות מהדאטהסט החשובים למשימה במהלך האימון של הרשת.

פרטי מאמר:

לינק למאמר: [זמין להורדה](#).

לינק לקוד: לא הצלחתי לאתר.

פורסם בתאריך: 5.11.20, בארקיב.

יוצג בכנס: ICLR 2021.

תחומי מאמר:

- שיטות אימון של רשתות נוירונים.
- שיטות מטה למידה (meta-learning) בתחום רשתות הנוירונים.

כלים מתמטיים, מושגים וסימונים:

- משפט הפונקציה הסתומה.
- חישוב נגזרת של פונקציה וקטורית דרך ההופכית של מטריצת הסיאן (hessian).
- קירוב ניומן (neumann) לחישוב הופכית של אופרטור (מטריצה) לינארי.
- רשת לומדת פנימית (inner student network).
- רשת מלמדת (נקראת הרשת המפרשנת במאמר - commentary network).
- אימון פנימי/חיצוני (inner/outer optimization).
- מטה-אימון, (meta-training).

תמצית מאמר:

כמו שאתם בטח יודעים, למרות הפעילות המחקרית האינטנסיבית בתחום הלמידה העמוקה עדיין קיימות לא מעט סוגיות פתוחות בנושאי אימון, רגולריזציה והבנה של מה שקורה בתוך רשתות נוירונים עמוקות. השאלות האלו נוגעות בסוגיות בסיסיות כמו: איך לאמן רשתות בצורה יותר מהירה, איך להקטין את כמות הדאטה הנדרשת לאימון, איך לשפר את יכולת הכללה רובסטיות של הרשתות.

אחת הגישות המעניינות שהוצעה לאחרונה שמנסה לתת מענה לשאלות אלו נקראת "לומדים ללמד" (learning to teach) המציע לבנות מנגנון חיצוני (רשת בדרך כלל) בשביל לספק לרשת הלומדת תובנות (נקרא לזה גם מטה-מידע בהמשך) לגבי המשימה תוך כדי תהליך האימון. למשל מנגנון כזה יכול לבצע משקול של דוגמאות במטרה לעזור לרשת הלומדת "לרכז את המאמץ" בדוגמאות החשובות. דוגמא אחרת של גישה זו יכולה להיות רשת עזר ה"מיעצת" איך לבנות דוגמאות (למשל ע"י ערבוב של דוגמאות מדאטה סט) הגורמים לרשת הלומדת לבנות ייצוג חזק של דאטה.

מאמר זה מציע מסגרת כללית לגישה זו (הנקראת למידה עם פרשנויות) ומציע תהליך אחיד להפקה של מטה-מידע (פרשנות) ע"י רשת חיצונית (מפרשנת) N_{com} , תוך כדי "הסקת מסקנות" העולות בתהליך האימון של רשת N_{st} (נקראת הרשת הלומדת) על סט האימון. אז בואו נבין איך כל זה עובד בעצם? נניח שאנו רוצים למצוא איזושהי טרנספורמציה (לדוגמא משקול/ערבוב) של דוגמאות בדאטה סט במטרה לשפר את הביצועים של הרשת הלומדת N_{st} והמטרה של הרשת המפרשנת N_{com} הינה למצוא את הטרנספורמציה הזו וזה למעשה מהווה הפלט שלה Out_{com} . במקרה זה תהליך האימון מכיל את השלבים הבאים:

- אופטימיזציה פנימית: עבור סט משקלים נתון W_{com} של הרשת N_{com} , מאמנים את N_{st} (כמה איטרציות של GD על משקלי W_{st}). במקרה הזה מפעילים טרנספורמציה Out_{com} המופקת ע"י (N_{com}, W_{com}) על הדאטה של סט האימון ומאמנים את N_{st} עליו. הפלט של השלב הזה הוא המשקלים W_{st} של N_{st} .
- אופטימיזציה חיצונית: מחשבים את הלוס של N_{st} עם סט המשקלים W_{st} מהשלב הקודם על סט ולידציה שעובר טרנספורמציה הניתנת ע"י " (N_{com}, W_{com}). כאן מאמנים את N_{st} על משקלי W_{com} כלומר מבצעים כמה איטרציות של GD אבל הפעם על למשקלי W_{com} .
- חוזרים על הצעדים אלו T פעמים כאשר T זה מספר האיטרציות של מטה אימון.

```

(1) Initialize commentary parameters  $\phi$  and student network parameters  $\theta$ 
(2) For  $M$  steps:
    (i) Compute the student network's training loss,  $\mathcal{L}_Y(\theta, \phi)$ .
    (ii) Compute the gradient of this loss w.r.t the student parameters  $\theta$ .
    (iii) Perform a single gradient descent update on the parameters to obtain  $\hat{\theta}$  (Note this is implicitly a function of  $\phi$ , i.e.  $\hat{\theta}(\phi)$ ).
    (iv) Compute the student network's validation loss,  $\mathcal{L}_V(\hat{\theta})$ .
    (v) Compute  $\frac{\partial \mathcal{L}_V}{\partial \phi}$ .
    (vi) Approximately compute  $\frac{\partial \mathcal{L}_Y}{\partial \phi}$  with equation 4, using a truncated Neumann series with a single term and implicit vector-Jacobian products [17].
    (vii) Compute the overall derivative  $\frac{\partial \mathcal{L}_V}{\partial \phi}$  using (v) and (vi), and update  $\phi$ .
    (viii) Set  $\theta \leftarrow \hat{\theta}$ .
(3) Output:  $\hat{\theta}$ , the optimized parameters of the commentary.

```

הסבר של רעיונות בסיסיים:

קודם כל נציין כי הגרדיאנט של משקלי N_{com} משערך את השינוי בלוס של N_{st} ביחס לשינוי במשקלים של N_{com} . אבל צריך לזכור שבשביל לחשב את הלוס של האופטימלי של N_{st} עבור משקלי N_{com} נתונים, המשקלים של N_{st} עוברים כמה איטרציות (אולי די הרבה) של GD במטרה למזער את הלוס שלה. לכן כדי לחשב את הגרדיאנט של הלוס של N_{st} לפי משקלי W_{com} צריך "לגלגל את כל האיטרציות על משקלי N_{st} " ל- W_{com} שזה יכול להיות די כבד חישובית כאשר N_{com} הינה רשת גדולה.

גם אם נחליט להשתמש רק באיטרציה אחת בתהליך האופטימיזציה הפנימית (זה מה שעשו במאמר) עדיין של לנו בעיה עם חישוב הגרדיאנט של הלוס לפי W_{com} . הבעיה הזו נובעת מהעובדה שגרדיאנט זה שווה למכפלה של הגרדיאנט של הלוס לפי W_{st} (שזה ניתן לחשב אותו בצורה הסטנדרטית של גזירת הלוס של רשתות) והנגזרת של וקטור משקלים W_{st} לפי לוקטור משקלים W_{com} . נזכיר ש W_{st} תלוי ב- W_{com} בצורה לא מפורשת כי בשלב האופטימיזציה הפנימית W_{st} מחושב על הדאטה סט אחרי הפעלת עליו טרנספורמציה המוגדרת ע"י W_{com} . נגזרת זו היא בעצם מטריצה (W_{com} ו- W_{st} הם וקטורים) שמימדיה עלולים להיות די גבוהים. נניח ש N_{st} ו- N_{com} הם רשתות לא גדולות בגודל של מיליון משקלים אז הנגזרת הזו תהיה מטריצה בגודל מיליון על מיליון ותידרש כמות זיכרון עצומה בשביל לאחסן אותה.

לכן המאמר מציע להשתמש במשפט הפונקציה הסתומה עבור הנגזרת של הלוס (של השלב החיצוני L_{out}) לפי W_{st} . משפט זה מאפשר לתאר את הנגזרת הבעייתית ע"י מכפלה של הופכית הסיאן של לוס של השלב הראשון L_{in} לפי W_{st} והמטריצה של הנגזרות המעורבות לפי W_{com} ו- W_{in} של L_{in} . למי שלא זוכר הסיאן זה מטריצה המורכבת מהנגזרות שניות של L_{in} לפי הזוג של רכיבים של W_{com} ו- W_{st} . צריך לזכור שהפירוק לעיל מתקיים בסביבת נקודה שבה הנגזרת של L_{in} לפי W_{st} מתאפסת. העובדה שלא ניתן למצוא אותה במדויק וזה יכול להשפיע בצורה שלילית על התהליך המטה-למידה.

גם אחרי הפירוק הזה יש לנו בעיה - והיא טמונה בחישוב של הופכית של הסיאן של L_{in} לפי W_{st} . אפילו עבור W_{st} בגודל יחסית לא גדול חישוב ההופכית (ולפעמים הסיאן עצמו) יכול להיות מאוד כבד וידרוש כמות עצומה של הזיכרון והזמן. בשביל להקל על ההיבט החישובי משתמשים ב**קירוב נוימן** עבור ההופכית מוכפלת בגרדיאנט של L_{in} לפי W_{st} (**מאמר של לוריין**) תוך שימוש בצורת עדכון של GD. דרך אגב נוסחת נוימן מאתרת הופכית של אופרטור לינארי בתור טור אינסופי של החזקות שלה (מוזוזות במינוס מטריצה היחידה). במאמר משתמשים רק באיבר אחד של קירוב זה.

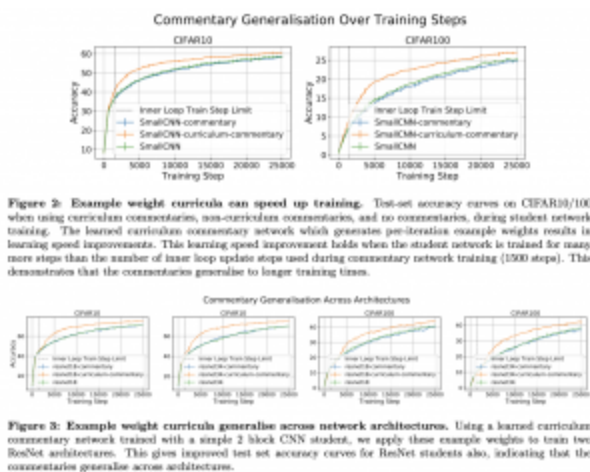
הישגי מאמר:

במאמר מראים 3 דרכים להשתמש בגישה זו לשיפור תהליך האימון של רשתות נוירונים:

- חישוב של משקול על דוגמאות מסט האימון ע"י W_{com} (דוגמאות עם משקל גבוה משפיעים יותר על הלוס). מעניין שהם גם בדקו את הביצועים של השיטה שלהם בתרחיש למידת few-shot על משימת מטח-למידה קלאסית. המטרה בלמידת few-shot היא לאמן רשת חיצונית (מטח) על מספר משימות (שנלמדות בפועל ע"י הרשת הפנימית) במטרה ללמד אותה להפיק תכונות משותפות של כל המשימות (שימו לב שזה מקרה פרטי של הפרדיגמה הכללית שהוצעה במאמר). כאשר מגיע משימה חדשה הרשת החיצונית מסוגלת לכייל את עצמה עם כמות קטנה של דאטה במשימה הזו. MAML זו אחת הדרכים לפתור בעיה זו והיא מאמנת רשת חיצונית לאתחול המשקלים של הרשת הפנימית שתאפשר לה להגיע לביצועים טובים על משימה חדשה במספר איטרציות GD קטן. אז הם מראים שהשילוב של MAML עם משקול דוגמאות הנבנה ע"י W_{com} גורם לשיפור ביצועים משמעותי. המאמר מראה שיפור בביצועים עבור הדאטה סטים MinilmaNet ו-CUB200-2011.

- בנייה של מקדמים ערבוב אופטימליים עבור הדוגמאות (בדומה ל mixup). כאן דוגמא מעורבת נבנית כסכום קמור של שתי דוגמאות: $x_{mix} = ax_1 + (1 - a)x_2$ כאשר מטרת W_{com} הינה לחשב את מקדמי a האופטימליים לביצועי משימת סיווג (מאמר mixup המקורי מגריל אותם מהתפלגות בטח). מעניין שכאן "למידה עם פרשנויות" מנצחת את mixup ב-CIFAR10 ומציגה ביצועים קצת פחות טובים ממנה על CIFAR10 (דאטה סט יותר קטן).

- חישוב מסכות על תמונות לשיפור הפיצ'רים המופקים ע"י רשת. כאן W_{com} בעצם מחפשת אזורים "חשובים" בתמונה שכדאי לרשת הלומדת להתרכז עליהם. המאמר מראה באופן ויזואלי שהמסכות שהוא מפיק אכן מתמקדות באזורים החשובים של תמונות ומראים באופן כמותי את עדיפותן על פני שיטות אחרות לבניית מסכות.



נ.ב. מאמר מציע מסגרת כללית לשיפור של תהליך למידה של רשתות נוירונים שניתן להשתמש בה למגוון רחב של משימות של הלמידה העמוקה. הגישה שלהם גם עוזרת להפיק תובנות חדשות תוך כדי תהליך האימון של רשתות. אני מני שעוד נשמע על שימושים רבים של גישה זו...

deepnightlearners#

הפוסט נכתב על ידי [מיכאל \(מייק\) ארליכסון](#), [PhD](#), Michael Erlihson.

מיכאל עובד בחברת סייבר [Salt Security](#) בתור Principal Data Scientist. מיכאל חוקר ופועל בתחום הלמידה העמוקה, ולצד זאת מרצה ומנגיש את החומרים המדעיים לקהל הרחב.