

לילה טוב חברים, היום אנחנו שוב בפינתנו deepnightlearners עם סקירה של מאמר בתחום הלמידה העמוקה
היום בחרתי לסקירה את המאמר שנקרא:
Unsupervised Learning of Visual Features by Contrasting Cluster Assignments
שיצא לפני בערך חודשיים

הוצג בכנס: NeurIPS 2020

תחומי מאמר:

- למידת ייצוג ללא דאטה סט מותיג (SSRL -self-supervised representation learning)
- SSRL מבוססת על טכניקות קליסטור (Clustering for deep representation learning)

כלים מתמטיים, מושגים וסימונים:

- מולטי-קרופ - טכניקת אוגמנטציה המבוססת על לקיחת פאטצ'ים קטנים של תמונה ברזולוציות נמוכות שונות
- האלגוריתם של סינקהורן קנופ (Sinkhorn-Knopp) לפתרון בעיית הטרנספורט האופטימלי למידות הסתברות דיסקרטיות

תמצית מאמר:

המאמר מציע שיטת למידת ייצוג על דאטה סט לא מותיג. רוב גישות המודרניות בתחום הזה (SSRL) מורכבות משני מרכיבים עיקריים:

הלוס המנוגד (contrastive loss - CL): מסתמך על ההנחה שייצוגים של דוגמאות קרובות צריכים להיות קרובים, בזמן שייצוגים של דוגמאות לא קשורות (נבחרות רנדומלית בד"כ) צריכים להיות רחוקים.
שיטה ליצירה של דוגמאות "דומות", קרי אוגמנטציה: בדרך כלל זוג דוגמאות קרובות (אקראי לזוגות האלו בהמשך זוגות חיוביים או זוגות קרו נוצרות ע"י ההפעלה של שתי אוגמנטציות שונות על אותה דוגמא.

נציין כי גישות SSRL המודרניות מסתמכות על השוואה של מספר גבוה מאוד של זוגות ייצוגים של דוגמאות שדורש כמות גדולה של זכרון ומשאבי עיבוד משמעותיים. דרישות אלו מקשות על יישום של שיטות אלו בצורה אונליין (לטענת המאמר רוב שיטות SSRL היום מיושמות בצורת אונליין שדי הפתיע אותי). אז בואו נדבר על החידושים שהמאמר הזה מציע:

שיטת אימון SwaV: המאמר הנסקר מציע שיטה חדשה SSLR (הנקראת SwaV) העשויה להוריד גם את כמות החישובים וגם לצמצם את כמות הזכרון הנדרשות. הרעיון העיקרי של המאמר הינו שינוי "ההגדרה של מושג הדמיון בין ייצוגי דוגמאות". למעשה המאמר "מאלץ" זוגות של הדוגמאות הקרובים "להשתייך" לאותם הקלאסטרים במרחב הייצוג במקום להשוות את הייצוגים בצורה מפורשת (שיוך זה המיוצג ע"י הקוד של דוגמא המחושב על סמך הבאטץ' שלו - אופן בנייתו יפורט בהמשך). נציין ש SwaV אינו דורש לשמור בנק של דוגמאות שליליות שהופך אותו למועמד טוב למימוש בצורה אונליין.

שיטת אוגמנטציה מולטי-קרופ: המאמר מציע שיטת אוגמנטציה הנקראת מולטי-קרופ שמתחילה מהחישוב של שני "קרופים סטנדרטיים" $1 \times \text{cr}$ ו- $2 \times \text{cr}$ של תמונה x . לאחר מכן לוקחים "קרופים קטנים יותר" של $1 \times \text{cr}$ ו- $2 \times \text{cr}$ במגוון רזולוציות נמוכות ובונים מהם סט דוגמאות חיוביות עבור תמונה x . לטענת המאמר שיטה זו מקטינה את כמות החישובים הנדרשת תוך שמירה על הביצועים

הסבר של רעיונות בסיסיים:

עכשיו ננסה להבין מה פונקציית המטרה L שבליבה של SwaV. פונקציית L מוגדרת באופן הבא (לכל דוגמא בבאטץ':)

- בונים מספר אוגמנטציות שלה לדוגמא x עם מולטי קרופ או כל גישה אחרת

- מרכיבים מאוגמנטציות אלו זוגות של דוגמאות
 - בונים וקטורי ייצוג z לכל הדוגמאות שבנינו
 - לכל זוג וקטורי ייצוג (z, z_1) מחשבים את הקודים שלהם 1_q ו- 2_q
 - מחשבים את סכום הדמיונות l_s בין 1_z ו- 2_q ובין 2_z ל 1_q
 - מחשבים את הסכום L_x של l_s של כל הזוגות של הדוגמאות החיוביות של דוגמא x
- אינטואיציה: למעשה תהליך אימון זה "מאלץ" וקטורי ייצוגי של דוגמא להכיל מידע על הקוד של הדוגמאות הקרובות. בצורה לא פורמלית ניתן לומר שאנו מנסים למקסם את "המידע הדדי" בין הייצוגים של הדוגמאות שזה המטרה העיקרית של האימון עם הלוס המוגד CL.
- דרך אגב השם של השיטה נובע מהפעולה שחלוף (swap) שמבצעים בין הייצוגים ובין הקודים של דוגמאות קרובות באימון.

השאלה האחרונה שטרם התייחסנו אליה הינה מבנה של פונקציית לוס בין ייצוג z לקוד q .
מבנה של פונקציה לוס בין ייצוג z לקוד q (של דוגמאות קרובות): אם אתם זוכרים הקוד q ניתן לפרש כווקטור הסתברויות שיוך לקלסטרים. למעשה אנו רוצים שהקוד q ישקף בצורה כמה שיותר טובה את המרחקים של z מהפרוטוטיפים c_j שניתן לראות אותם בתור מרכזים (סנטרואידים) של קלסטרים של ייצוגים. אז קודם כל אנו בונים את וקטור המרחקים המנורמלים מ- z לכל c_j . מרחק זה מחושב כאקספוננט של המכפלה הפנימית בין z ל c_j . בסוף לוקחים את וקטור המרחקים ומנרמלים אותו. לאחר מכן מחשבים את קרוס אנטרופי בין q לווקטור מרחקים מנורמל שחישבנו. את הפונקציה זו אנו ממקסמים ביחס ל ייצוגים z וביחס לפרוטוטיפים c .

אינטואיציה: שימו לב על הדמיון של המרחק בין וקטור הייצוג z ל- c_j לביטוי של החוב המוגד CL. וזה לא מקרי - אתם זוכרים שלהבדיל משיטות מבוססות CL קלאסי, אין לנו כאן דוגמאות שליליות בצורה מפורשת. אז מה שמשחק כאן את תפקיד "הדוגמאות השליליות" זה מרכזי הקלסטרים שרחוקים מ z . כלומר הם מאלצים ייצוגים של דוגמאות חיוביות להיות רחוקים בצורה כמה שיותר דומה מכל הקלסטרים השליליים וקרובים באותה מידה מהקלסטרים החיוביים. לדעתי זה הנקודה הכי משמעותית במאמר (!!).

הסבר על בניית קוד q של ייצוג z : הקוד q של וקטור ייצוג z מתאר את "רמת קרבתו" של z ל K וקטורי פרוטוטיפ c_j . וקטור c_j "מייצג" את הקלסטר i . קוד של דוגמא (וגם של כל האוגמנטציות שלה) מחושב על סמך הבאטץ' בלבד (!!). אפשר להגיד שהקוד q מייצג את ההסתברויות שיוך של וקטור הייצוג z של הדוגמא לקלסטרים המיוצגים ע"י וקטורי c_k .

מטריצה Q המכילה את הקודים של כל הדוגמאות מהבאטץ' הינה פתרון של בעיית אופטימיזציה לינארית עם איבר רגולריזציה השווה לאנטרופיה הכוללת של Q עם מקדם קטן. פונקציה מטרה זו מנסה למקסם את הדמיון הכולל בין וקטורי ייצוג של הדוגמאות בבאטץ' לפרוטוטיפים c_j (כלומר לפזר את הקודים בצורה המשקפת את את יחס המרחקים בין ייצוג הדוגמא למרכזי הקלאסטרים השונים). שימו לב שבעיית אופטימיזציה זו מזכירה בצורתה את בעיית הטרנספורט האופטימלי בין מידות הסתברות דיסקרטיות (האחידות) המוגדרות על שני דאטה סטים. את התפקיד של דאטה סטים כאן משחקים הפרוטוטיפים c_j וקטורי הייצוג z_j של כל הדוגמאות בבאטץ'. המטרה כאן זה למצוא את האופן האופטימלי שבו ניתן "להעביר את המסה ההסתברותית מווקטורי z_j לוקטורי c_j (נציין שפונקציית המרחק שיש בהגדרה של הטרנספורט האופטימלי הינה פרופורציונלית במקרה שלנו למרחק בין z ל c). למעשה אנו מנסים למצוא מטריצה Q האי שלילית, שאיבר jk שלה מגדיר את המסה ההסתברותית המועברת מווקטור z_k לווקטור c_j כלומר הסתברות השיוך של z_k לקלסטר של c_j . מכיוון שאנו רוצים שאותו מספר דוגמאות "ישוין" לכל קלאסטר, מוסיפים אילוץ על סכום השורות וסכום העמודות של Q . בעיה זו פותרים בעזרת אלגוריתם איטרטיבי של סינקהורן-קנופ.

הסבר על מושגים חשובים במאמר:

שיטות אימון של גישות SSRL המודרניות: בדרך כלל בזמן האימון של SSLR לכל זוג של דוגמאות קרובות בונים מספר גדול של זוגות רנדומליים (אקרא לזוגות כאלו זוגות רחוקים או זוגות שליליים). כאן פונקציית המטרה F_{ob}

(שממקסמים אותה) הינה יחס בין אקספוננט של דמיון של "הזוג הקרוב" (בין הייצוגים שלהם) לסכום הדמיונות בינו לבין כל הזוגות שליליים. למשל שיטת [SimCLR](#) כל באטץ' מורכב מ-N זוגות של דוגמאות קרובות (אוגמנטציה של אותה הדוגמא) המהווים את הזוגות החיוביים כאשר עבור דוגמא נתונה, כל הדוגמאות פרט ל"בת הזוג" שלה נחשבת לדוגמא שלילית עבודה. פונקציה המטרה לכל באטץ' הינה סכום של פונקציות המטרה של כל N2 דוגמאות של הבאטץ'.

בנק של ייצוגי דוגמאות שליליות: ידוע שהגדלת מספר הזוגות השליליים לכל זוג חיובי באימון תורמת לעוצמת הייצוג של הדאטה. כתוצאה מכך משתמשים בבאטצ'ם מאוד גדולים (עשרות אלפי דוגמאות) שדורש משאבי זכרון גדולים, כח עיבוד רב (צריך לחשב את הייצוג של עשרות אלפי דוגמאות מהבאטץ'). כדי להקטין את כוח העיבוד הנדרש הוצע ([MOCO](#)) "בנק הדוגמאות השליליות" מהבאטצ'ים הקודמים המכיל את הייצוגים של הדוגמאות מכמה הבאטצ'ים הקודמים. כל פעם דוגמים משם ייצוגים של דוגמאות שליליות ומוספים את זה לייצוגיים השליליים מהבאטץ' הנוכחי. צריך לזכור שגישה זו כרוכה בהקצאת משאבי אחסון נוספים לשמירת בנק זה.

הישיגי מאמר: המאמר מראה ש SwaV משולב עם מולטי-קרופ מצליח לייצר ייצוגים יותר חזקים משיטות בניית ייצוג רבות במספר משימות. ההשוואה בוצעה בדרך הסטנדרטית: הוספה של שכבה לינארית לרשת הבונה ייצוג(עם משקלים מוקפאים) ובחינת ביצועיה על משימה מסוימת. קודם כל הם הראה שייצוג שנבנה באמצעות SwaV מציג ביצועים יותר טובים על דאטה סטים 2018VOC07 iNaturalist ו- 205Places מהייצוגים הנבנים על ImageNet מתווג (!!) גם על משימת סיווג ועל משימת זיהוי אובייקטים. בנוסף הם הראו שהייצוגים שלהם משיגים ביצועים יותר טובים מבחינת 5Top1/Top (לוקחים 5/1 דוגמאות הכי קרובות מבחינת הייצוג ומחשבים כמה מתוכם שייכים לאותה קטגוריה) משיטות כמו SimCLR ו- 2MoCov. נזכיר שלהבדיל מ 2MoCov, אין צורך בשמירה של בנק דוגמאות שליליות ב SwaV. הם גם הראה את עליונותה של SwaV במשימות semi-supervised על שיטות כמו UDA ו- FixMatch. וזה רק חלק קטן מכל השוואות שהם עשו - הם באמת עשו עבודה מרשימה בהיבט הזה.

לינק למאמר: <https://arxiv.org/pdf/2006.09882.pdf>

לינק לקוד: <https://github.com/facebookresearch/swav>

נ.ב. מאמר ממש מגניב עם רעיון מגניב המשלב תובנות רבות ממגוון שיטת SSRL. הם גם טרחו להשוות את הביצועים של השיטה שלהם מול מגוון רחב של אלגוריתמים, משימות, דאטה סטים וקונפיגורציות שזה בהחלט מרשים. בקיצור המלצת קריאה לווהטת ממני:

#deepnightlearners