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:

AVAE: Adversarial Variational AutoEncoder

## : 21.12.2020 תאריך פרסום

הוצג בכנס: טרם ידוע

## :תחומי מאמר

- (VAE Variational AutoEncoder). אוטו-אנקודר וריאציוני
  - .(GANs Generative Adversarial Networks) גאנים •

## כלים מתמטיים, טכניקות, מושגים וסימונים:

- פונקצית לוס של VAE (המתקבלת מ- VAE) (המתקבלת של VAE)
  - מרחק KL בין התפלגויות.
  - Mutual Information).) מידע הדדי בין משתנים אקראיים/התפלגויות
    - information bottleneck).) צוואר בקבוק אינפורמציוני
- פונקציית הלוס הסטנדרטית של גאן (מ<u>המאמר המקורי</u>) והפתרון האופטימלי שלה מבחינת הדיסקרימימטור.

בהירות כתיבה: בינונית מינוס

רמת היכרות עם כלים מתמטיים וטכניקות של ML/DL הנדרשים להבנת מאמר: נדרשת הבנה עמוקה ב- VAE, גאנים ותכונותיהם בשביל להבין לעומק את הרעיון הבסיסי של המאמר. שליטה בכלים מתמטיים מתחום ההסתברות והסטטיסטיקה נחוצה להבנת המאמר.

יישומים פרקטיים אפשריים: גינרוט תמונות באיכות גבוהה עם VAE (סוג של 😛 ).

המלצת קריאת ממייק: מומלץ לבעלי ידע עמוק ב- VAE, גאנים ובעלי ידע מוצק בהסתברות בתור אתגר.

מבוא והסבר כללי על תחום המאמר: יצירה של תמונות פוטוריאליסטיות עי" רשתות נוירונים הפכה לנושא חם בלמידה העמיקה מאז שיאן גודפלו (Ian Goodfellow) הגה את הרעיון של 2018. לנושא חם בלמידה העמיקה מאז שיאן גודפלו לצירת דאטה במספר דומיינים שהפופולריים 2014. מאז הוצעו מספר מודלים גנרטיביים שונים ליצירת דאטה במספר דומיינים שהפופולריין ביצירת CAN, כאשר לכל אחד מהם יתרונות וחסרונות משלה. למשל מאין ביצירת תמונות שנראות ממש כמו אמיתיות (קרי פוטוריאליסטיות) אך הוא מאוד קשה לאימון. התופעות כמו Mode Collapse (יצירה של תמונות כמעט זהות עי" הגנרטור) וגם ההתכנסות של תהליך האימון אינה מובטחת – אלו רק חלק מהבעיות שעלולות לעלות במהלך אימון של גאן. בנוסף המבנה של המרחב הלטנטי של גאן אינו נוח לניתוח ולא נתון בצורה מפורשת. מהעבר שני VAEs יותר קלים לאימון והמרחב הלטנטי שלהם נתון בצורה מפורשת יותר אך התמונות שנוצרות באמצעותם הן מטושטשות ופחות פוטוריאליסטיות לרוב.

קודם כל ניזכר ממש בקצרה מה זה בעצם VAE.

הסבר קצר על :VAE ארכיטקטורה של VAE מורכבת משתי רשתות עם פרמטרים נלמדים.

- Z שממפה דוגמא מהמרחב המקורי למרחב הלטנטי E\_vae (אנקודר)(בעל מימד נמוך).
  - הרשת המפענחת D\_vae (דקודר) מנסה לשחזר את הדוגמא מהייצוג הלטנטי שלה.

האימון של VAE מתבצע בצורה הבאה:

- .z ממפה דוגמא X לפרמטרים של הייצוג הלטנטי שלה E vae הרשת המקודדת
  - . מגרילים וקטור אקראי (בדרך כלל גאוסי) עם הפרמטרים מהשלב הקודם.
- שחזור של הדוגמא המקורית D\_vae מעבירים את הוקטור המוגרל דרך הרשת המפענחת ∆ .X

פונקצית הלוס של VAE, המסומנת עי" L\_vae, מתקבלת עי" שימוש בחסם העליון של – ELBO – פונקצית הלוס של sevidence. פונקצית לוס זו מורכבת משני מחוברים:

- לוס השחזור: עד כמה טוב D\_vae הצליחה לשחזר את התמונה המקורית X. בדרך כלל לוס
  השחזור מחושב כמרחק הריבועי L\_q בין התמונה בין התמונה המקורית למשוחזרת.
- מרחק KL בין התפלגות פריור על מרחב הלטנטי לבין התפלגות פוסטריור שלה (המשוערכת על סמך פלטים של הרשת המקודדת E vae).

בתהליך האימון של VAE הרשת המקודדת והרשת המפענחת מאומנות במטרה למזער את L vae

כדי להבין את הסיבות העיקריות ליכולת החלשה של VAE ליצור תמונות פוטוריאליסטיות, המאמר מנתח את פונקצית הלוס שלה ומציין שתי סיבות עיקריות לכך:

<u>סיבה 1:</u> המאמר מנתח את האיבר השני שלו, כלומר L\_kl, ומוכיח שניתן לתאר את L\_kl כסכום של המידע ההדדי בין הדוגמה x לייצוג הלטנטי האפוסטריורי שלה z|x, ומרחק KL בין ההתפלגות האפוסטריורית z|x והתפלגות הפריור של z (בדרך כלל גאוסי בעל תוחלת אפס מטריצת קווריאנס (l. מזעור של איבר זה משמעו

לזו של z|x במטרה לקרב את ההתפלגות של z|x הגבלה על מידע הדדי בין דוגמא לייצוגה הלטנטי (!!) במטרה לקרב את ההתפלגות של הפריור z.

כלומר, איבר זה הינו למעשה צוואר בקבוק אינפורמציוני המגביל את זרימת המידע בין x לייצוג הלטנטי שלה. זה בפועל מקשה על הרשת המפענחת לשחזר את התמונה המקורית מהקוד הלטנטי שלה (כי חלק מהמידע הולך לאיבוד בין הדוגמא לייצוגה הלטנטי עקב צוואר הבקבוק האינפורמציוני). בנוסף הלוס הריבועי המופיע באיבר הראשון של L\_vae, גורמת ל- VAE ליצור תמונות מטושטשות. המאמר מצטט עבודה של ECun et al המראה שלמעשה הערך האופטימלי של כל פיקסל בתמונה המשוחזרת הינו לו התוחלת שלו המותנית ב"מידע הנמצא בקוד הלטנטי שלה". כתוצאה מכך הרשת המפענחת לרוב פשוט לא מצליחה להפיק תמונה פוטוריאליסטית מהייצוג החלקי שמוזן אליה.

סיבה 2: הסיבה השנייה טמונה בהנחה המקובלת בראייה הממוחשבת כבר שנים: לתמונות הטבעיות low- יתירות רבה (במיוחד היתירויות הלוקאליות) שמאפשרת לתאר אותן במרחב ממימד נמוך (-low- יתירות רבה (במיוחד היתירויות הלוקאליות) שמאמר) הטקסטורות כמו עץ, שער או גלים אינם "חיים" במניפולד ממימד נמוך, עקב התכונות האינהרנטיות שלהם (למשל שיער של אדם מורכב מהרבה שערות שלכל אחד אופיינים משלו). לעומת VAE, הגאנים מצליחים להתגבר חלקית על סוגיה זו ע"י יצירה של דוגמאות המהוות תת-קבוצה של המניפולד ממימד גבוה שבו "חיות" הטקסטורות, שמספיקות כדי "להוליך שולל" את הדיסקרימינטור. כתוצאה מכך התמונות של גאן יוצאות יותר "טבעיות" ופחות מטושטשות של VAEs.

רעיון המאמר בגדול: המאמר מציע לשלב את היתרונות של VAE וגאנים עי" שיבוצם בארכיטקטורה שלהם, הנקראת AVAE. אני רוצה לציין שהרעיון הזה לא חדש וכבר ב- 2015 ניסו לעשות זאת ב- VAE/GAN. במאמר הציעו להחליף את השגיאה הריבועית ברשת המפענחת בדיסקרימינטור כמו זה של גאן. הבעיה בגישה הזו שהדמיון של התמונה המשוחזרת למקורית כבר לא בא לידי ביטוי. שיטה של גאן. הבעיה בגישה הזו שהדמיון של התמונה המשוחזרת למקורית כבר לא בא לידי ביטוי. שיטה יותר מפורסמת המשלבת את שתי גישות אלו (VAE) חלקית וגאן) הינה BiGAN שמורכב מהרשת המקודדת, הרשת המפענחת והדיסקרימינטור. הדיסקרימינטור מנסה להבחין לא רק בין התמונות המגונרטות לאמיתיות אלא גם בין הווקטורים הלטנטיים שנדגמים המפריור לבין לאלו שנוצרים עי" הרשת המקודדת. BiGAN מצליח ליצור תמונות פוטוריאליסטיות אך (לטענת המאמר) יכולת השחזור

שלו נמוכה (כלומר תמונות נוצרות מוקטורים קרובים לייצוג לטנטי של תמונה נתונה x, לא תמיד יוצאות דומות ל- x), כלומר המרחב הלטנטי פחות קוהרנטי.

המאמר הנסקר מציע לשלב את לוס השחזור עם הלוס האדוורסרי בצורה שתיצור גם תמונות באיכות דומה בלי לפגוע ב״קוהרנטיות של המרחב הלטנטי״. הרעיון העיקרי של AVAE הינו הוספת רשת דומה בלי לפגוע ב״קוהרנטיות של המרחב הלטנטי״. הרעיון העיקרי של VAE המקודדת הגנרטור G ל- VAE הסטנדרטי, שלוקחת כקלט את הווקטור הלטנטי המופק עי״ הרשת התפלגות E\_vae המאמר גם מציע להוסיף(zoncatenate) לקלט של G וקטור נוסף בענת המאמר ב z\_a מיועד לייצוג מידע מתמונה שהרשת E\_vae לא הצליחה להפיק ממנה.

## הסבר מעמיק על רעיונות בסיסיים:

נתחיל מההסבר על פונקצית לום L\_avae של L\_avae פונקצית הלוס של AVAE הינה סכום L\_vae פונקצית הלוס של AVAE המסומן ב- L\_Vae. והלוס של הגנרטור G, המסומן ב- L\_vae עקרונית הגנרטור G צפוי לשרת כ"הופכית" של הרשת המקודדת E\_vae וזה הנקודה החשובה של המאמר:

ההתפלגות המותניתZ (p\_G(x|z) של פלט הגנרטור G בהינתן וקטור לטנטי z והתפלגות (p\_enc(x|z ההתפלגות g בהינתן וקטור לטנטי z, צריכות להיות כמה שיותר קרובות (כאשר של פלט(!!) הרשת המקודדת, בהינתן וקטור הלטנטי z,, צריכות להיות כמה שיותר קרובות (כאשר הווקטור z מתפלג לפי התפלגות הפריור).

(p\_enc(x|z לומר פונקצית המטרה G של הגנרטור G של הגנרטור L\_G של הערופי בין עלומר פונקצית המטרה G של הגנרטור G של ההתפלגות בין G מתפלג לפי התפלגות (log(p\_enc(x|z של התפלגות ב'p\_G(x|z של התפלגות D(0, I)).

הערת צד לגבי האימון :במהלך האימון של AVAE, הקלט של G אינו נלקח מהפלט של הרשת המקודדת E\_vae , אלא נדגם ישירות מהתפלגות הפריור של z. אני מנחש שזה הופך את הגנרטור יותר דומה לזה של הגאן המקורי במטרה "לחקות" את תכונות החזקות שלו ביצירה פוטוריאליסטיות.

נתחיל מקצת אינטואיציה מאחורי הרעיון הדי לא טריוויאלי הזה:

פינת אינטואיציה: שימו לב שהמטרה כאן היא לאמן את הגנרטור ליצור תמונות פוטוריאליסטיות מהתפלגות הפריור מחד (נראה בהמשך איך LG מוביל ללוס דומה לזה של הגאן הסטנדרטי) עבור וקטורים לטנטיים המתפלגים לפי הפריור הנתון. שנית זה "מאלץ" את הרשת המקודדת להפיק פיצ'רים לטנטיים הנחוצים ליצירת תמונה פוטוראליסטית. שלישית L\_vae מזעור המרחק הריבועי בין התמונה E\_enc להפיק פיצ'רים הנחוצים לשחזור מדויק של התמונה (עי" מזעור המרחק הריבועי בין התמונה המקורית למשוחזרת). המשחק המורכב הזה מאפשר ל- G ליצור תמונות פוטוריאליסטיות מחד תוך שמירת קוהרנטיות של המרחב הלטנטי (וקטורים לטנטיים קרובים יוצרים תמונות דומות כלומר "יחסי המרחק" במרחב המקורי ובמרחב הלטנטי נשמרים).

עכשיו בואו נבין איך L\_G מוביל ללוס דומה לזה של גאנים ״הגורם״ לו ליצור תמונות פוטוריאליסטיות. המאמר מוכיח ש- L\_G ניתן לפירוק לסכום של המחוברים הבאים:

- האיבר הראשון: תוחלת של (log(p\_enc(z|x\_G)) מעל ההתפלגות (yog(p\_enc(z|x\_G)) ב האיבר הראשון: תוחלת של (log(p\_enc(z|x\_G)) מעל ההתפלגות הפריור. איבר זה למעשה משערך עד כמה "סביר" להפיק וקטור לטנטי z מהתמונה דרך הרשת המקודדת E\_vae, כאשר x\_G שנוצר עי" הגנרטור בהינתן אותו וקטור לטנטי z. החלק הזה הוא קל יחסית לחישוב ומשוערך כמרחק ריבועי בין z לבין אותו וקטור לטנטי z. החלק הזה הוא קל יחסית לחישוב ומשוערך כמרחק ריבועי בין z לבין בהעותו של הפלט שלו. נזכיר שבהינתן התוחלת של הפלט ש- zog מפיק מ- x\_G מנורמל בשונות של המקודד מוציא זוג של (mu\_enc(x), sigma\_enc(x)) של התוחלת והשונות של הייצוג הלטנטי של z בהתאמה (מהם מגרילים את הקלט לרשת המפענחת zog).
- האיבר השלישי: אנטרופיה של p\_G(x) (התפלגות התמונות הנוצרות עי" הגנרטור). איבר זה אינו תלוי בדאטה (מתאר את ההתפלגות הפלט הגנרטור). מינימיזציה של איבר זה מאלצת את מתאר את הרוכזת סביב המודים שלה (זה גורם לירידה באנטרופיה) והמאמר (intractable מציע לא לקחת אותו חשבון כאמצעי רגולריזציה (נציין שהאיבר הזה הינו

לסיכום, AVAE מורכב מ 4 הרשתות הבאות:

- הרשת המקודדת הסטנדרטית E vae
- .D vae הרשת המפענחת הסטנדרטית
- הגנרטור G המיועד ליצירת תמונה מוקטור לטנטי (עם אופציה להוסיף וקטור לטנטי נוסף G הגנרטור של תמונות ש- E\_vae לכיסוי של תכונות של תמונות ש-
- הרשת המבקרת (דיסקרימינטור) C שמטרתה להבחין בין התמונות מהדאטהסט לתמונות מגונרטות

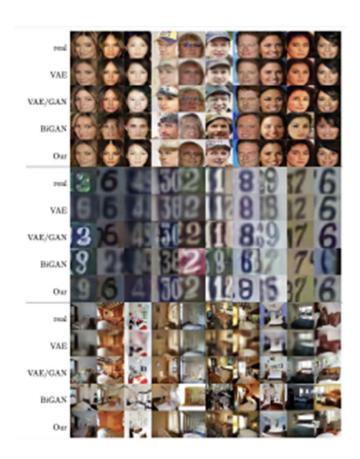
המפלצת הזו מאומנת עם פונקציית מטרה שהיא סכום של L\_G ו- L\_vae שמוסברים מעלה.

VAE, VAE/GAN, מול השיטות הבאות: המאמר משווה את ביצועיו של AVAE מול השיטות הבאות: המאמר משווה את ביצועיו של BiGAN ומשתמש ב 3 מטריקות: לוס השיחזור הריבועי, BiGAN

- Bedroom
  - <u>CelebA</u> ●
- CIFAR10 •
- CIFAR100
  - SVHN •

עבור כמה מהדאטהסטים האלו הם הצליחו להשתפר בחלק מהמטריקות (הלוס הריבועי שופר עבור כל הדאטהסטים) אבל רוב השיפורים לא נראים לי ממש מרשימים. הם גם טוענים שהתמונות שלהם כל הדאטהסטים) אבל רוב השיפורים לא נראים ואני נוטה להסכים איתם אך זה די סובייקטיבי נראות יותר פוטוראליסטיים מאלו של VAEs





לינק למאמר: זמין להורדה

לינק לקוד: למרות שבמאמר מופיע שהקוד יהיה זמין בגיטהאב לא הצלחתי לאתרו.

נ.ב. מאמר עם רעיון מאוד מעניין ומגניב אך כתוב בצורה לא מספיק ברורה. תרשימי זרימה שיש במאמר לרוב לא עוזרים בהבנה. צריך להודות שהתוצאות לא הרשימו אותי יותר מדי וגם הקוד לא שותף שזה די מאכזב. נאלץ לא להעניק לו המלצת קריאה ממייק. עם זאת הכלים/תובנות המתמטיים שפותחו במאמר נראים לי מבטיחים ומעניינים ואני מקווה לראותם מיושמים בעתיד ומשיגים תוצאות משמעותיות יותר.

#deepnightlearners