

רשתות נוירונים לתמונות תרגיל 3

יואב שפירא 312492838

23 במאי 2023

חלק מעשי

שאלה 1 *Loss Saturation*

ארכיטקטורות:

- *Generator*: מורכב מ-4 שכבות של *Strided Transpose Convolution* עם קרנל בגודל 3. אחרי כל שכבה מה-3 השכבות הראשונות יש שכבת *BatchNormalization* ואקטיבציה *ReLU*. לאחר השכבה הרביעית אין *BatchNormalization* והאקטיבציה היא *Sigmoid*.
בשכבה הראשונה, האינפוט (וקטור בגודל 100) ממופה ל-256 ערוצים לתמונה בגודל 4×4 .
בשכבה השנייה ממופה ל-128 ערוצים, בשלישית ל-64 ערוצים וברביעית לערוץ 1.

```
Generator(  
  (seq): Sequential(  
    (0): ConvTranspose2d(100, 256, kernel_size=(3, 3), stride=(2, 2), bias=False)  
    (1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)  
    (2): ReLU(inplace=True)  
    (3): ConvTranspose2d(256, 128, kernel_size=(3, 3), stride=(2, 2), bias=False)  
    (4): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)  
    (5): ReLU(inplace=True)  
    (6): ConvTranspose2d(128, 64, kernel_size=(3, 3), stride=(2, 2), bias=False)  
    (7): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)  
    (8): ReLU(inplace=True)  
    (9): ConvTranspose2d(64, 1, kernel_size=(3, 3), stride=(2, 2), padding=(2, 2), output_padding=(1, 1), bias=True)  
    (10): Sigmoid()  
  )  
)
```

- *Discriminator*: מורכב מ-3 שכבות של *Strided Convolution* עם קרנל בגודל 4, ולבסוף שכבת *FC* שממפה לאאוטפוט - מספר יחיד שהוא ההתסברות שהתמונה היא אמיתית. אחרי כל שכבת קונבולוציה יש *BatchNormalization* ואקטיבציה *LeakyReLU* עם שיפוע של 0.2 לערכים שליליים. אחרי שכבת *FC* יש שכבת אקטיבציה *Sigmoid* כדי למפות את הערך האחרון להסתברות.
בשכבה הראשונה האינפוט (תמונה בגודל 28×28) ממופה ל-32 ערוצים, בשכבה השנייה ל-64 ערוצים, ובשכבה השלישית ל-128 ערוצים. לבסוף ה-*FC* ממפה לערך יחיד.

```

Discriminator(
  (seq): Sequential(
    (0): Conv2d(1, 32, kernel_size=(4, 4), stride=(2, 2), padding=(1, 1))
    (1): BatchNorm2d(32, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (2): LeakyReLU(negative_slope=0.2, inplace=True)
    (3): Conv2d(32, 64, kernel_size=(4, 4), stride=(2, 2), padding=(1, 1))
    (4): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (5): LeakyReLU(negative_slope=0.2, inplace=True)
    (6): Conv2d(64, 128, kernel_size=(4, 4), stride=(2, 2), padding=(1, 1))
    (7): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (8): LeakyReLU(negative_slope=0.2, inplace=True)
    (9): FlattenBatch()
    (10): Linear(in_features=1152, out_features=1, bias=True)
    (11): Sigmoid()
  )
)

```

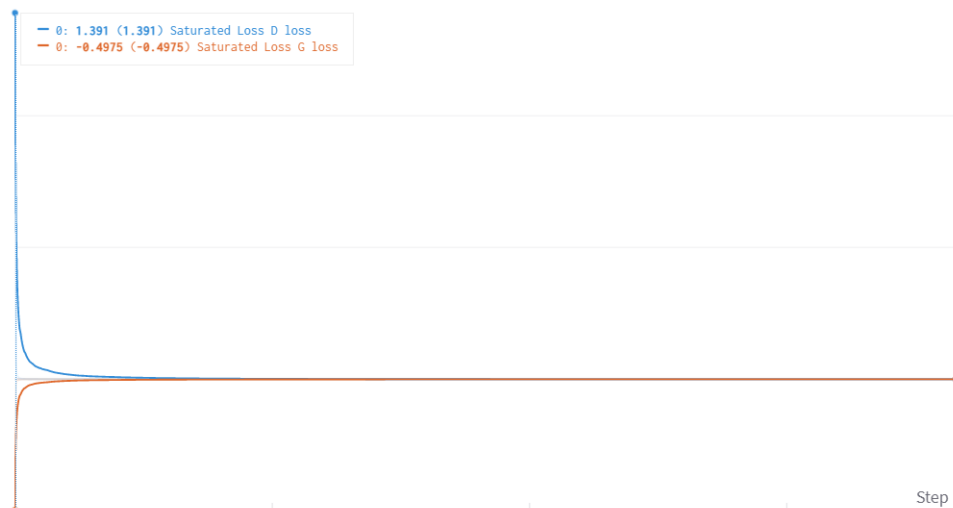
הבחירה לבנות את D קטן יותר מ G נבעה מניסוי ותעייה: עבור D גדול יותר, המודל הגיע למצב *Mode Collapse* לאחר מספר אפוקים של אימון. החלשה של D גרמה לו להתכנס לאט יותר, וכך אפשרה ללמידה להמשיך לאורך יותר זמן ללא *Mode Collapse*.

אימון:

השתמשתי באופטימיזר *Adam* עבור שני המודלים, עם רגולריזציה (נורמת l_2) עם ערך $\lambda = 0.5$. אימנתי ב-10 אפוקים, עם $lr = 0.002$. על כל 3 באצ'ים של אימון ל D אימנתי עם באצ' אחד את G . הדאטא עליו אימנתי את המודלים הוא כל הדאטאסט של *MNIST*, גם הטריין וגם הטסט.

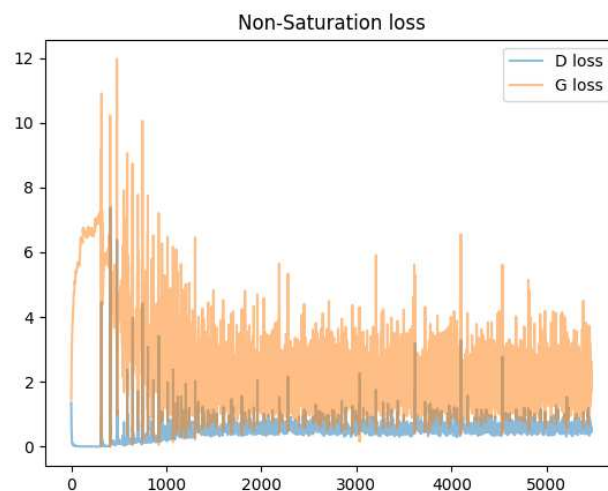
שאלה 1 פונקציות Loss:

1. *Original GAN Cross – Entropy*: בשביל ה $loss$ הזה השתמשתי ב= $criterion$ *BCELoss* כפונקציית $loss$ של D . כעת, G צריך למזער את השגיאה של $\log(1 - D(G(z)))$. כדי להשיג את זה השתמשתי באותה *BCELoss* עם הפרדיקציות של D על התמונות המזויפות, ועם $response$ של "תמונות מזויפות", ומכיוון שזה *Binary Cross Entropy* (מחושב על שני ערכים בלבד) נשיג את האפקט הרצוי. הגרף של $Losses$ של שני המודלים מוצג כאן, בכתום G ובכחול D (הערכים במלבן הם ההתחלתיים, ואפשר לראות שהגרדיאנטים מתאפסים מהר):



התוצאות שהמודל הזה הניב הן תוצאות של רעש: G לא הצליח לייצר תמונות כי בכלל לא היה לא זמן ללמוד. לא צרפתי תמונה זה לא מעניין זה פשוט רעש.

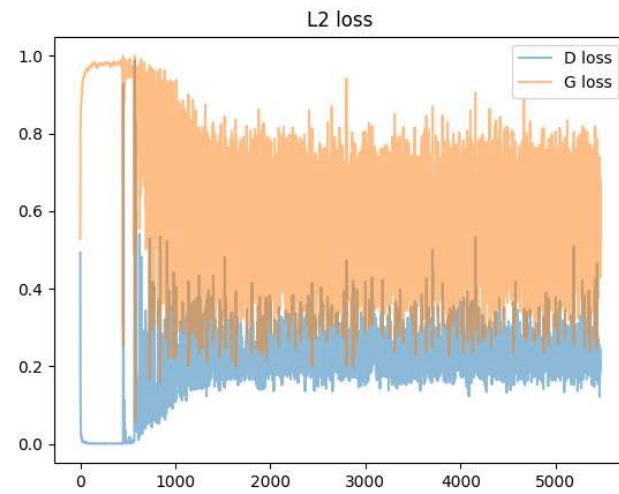
2. *Non-Saturation Loss*: גם כאן השתמשתי ב*BCELoss* כ*criterion*, רק שהפעם G צריך למקסם את $\log(D(G(z)))$, כלומר למזער את $-\log(D(G(z)))$, וכאן האפקט מתקבל על ידי מילוי הוקטור *response* בלייבלים של "תמונות אמיתיות". בגרפים של *Losses* אפשר לראות שיש התכנסות למצב מסוים - המצב של *Nash Equilibrium* - המודל D מאוד מוצלח בהתחלה, ולאט לאט הוא שוגה, והמודל G בדיוק הפוך, עד ששניהם יחסית מתאזנים כלומר מצליחים 'לבלבל אחד את השני'. בשיווי משקל אולטימטיבי, הדיוק של D יהיה 0.5, וניתן לראות שהוא מתקרב לשם.



התוצאות של המודל: לאחר 10 אפוקים G הצליח לייצר תמונות שנראות די טוב:



3. $L2 Loss$: השתמשתי ב $MSELoss$, ובחישוב ה $loss$ של G שלחתי את הפרדיקציות של D על תמונות מזויפות מ G , עם לייבלים של "תמונות אמיתיות". בגרף הזה אפילו יותר קל לראות את הכיוון לשיווי משקל נאש, כי השגיאה של שני המודלים חסומה בין 0 ל:1:



התוצאות לא רעות:



פונקציית $Loss$ שהגיעה לסטורציה היא הפונקצייה הראשונה - המקורית. זה קורה בגלל שהנגזרת במקרה הזה $-\frac{1}{1-D(G(z))}$, צריך לשים לב שבשלבם ההתחלתיים G מייצר רק תמונות מזויפות - ולכן D קל להבין את זה מאוד מהר. לכן, הגרדיאנטים בהתחלה יהיו מאוד קרובים ל-1, לעומת מצב של $Non - Saturation$ ובו הנגזרת שווה ל- $-\frac{1}{D(G(z))}$, ולכן הגרדיאנטים בהתחלה יהיו קרובים ל- $-\infty$. מכיוון ש- $-\infty$ הרבה יותר רחוק מ-0 מאשר 1, אנחנו מקבלים שעבור $loss$ הראשון הגרדיאנטים נעלמים מאוד מהר לעומת $loss$ השני. המצב של סטורציה הוא באמת מצב של $Vanishing Gradients$ בשלב התחלתי באימון של G .

שאלה 2 Model Inversion:

בשאלה זו לקחתי את המודל G שאומן עם $Non - Saturation Loss$, והקפאתי אותו בעזרת $eval()$ כדי שלא יתעדכן במהלך האופטימיזציה. הגרלתי z וקטור רעש, והגדרתי

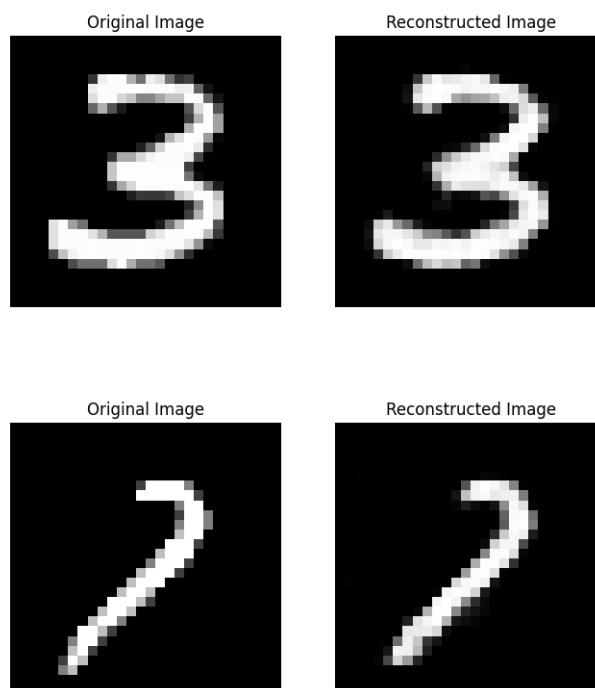
$Optimizer = Adam$ עם הפרמטרים של z . בעזרת $criterion = MSELoss$ חישבתי את $Loss$ של $|G(z) - I|$ עבור תמונה אקראית I מתוך הדאטאסט. לאחר 1000 איטרציות של GD קיבלתי z שמייצר את התמונה המקורית I יחסית טוב.

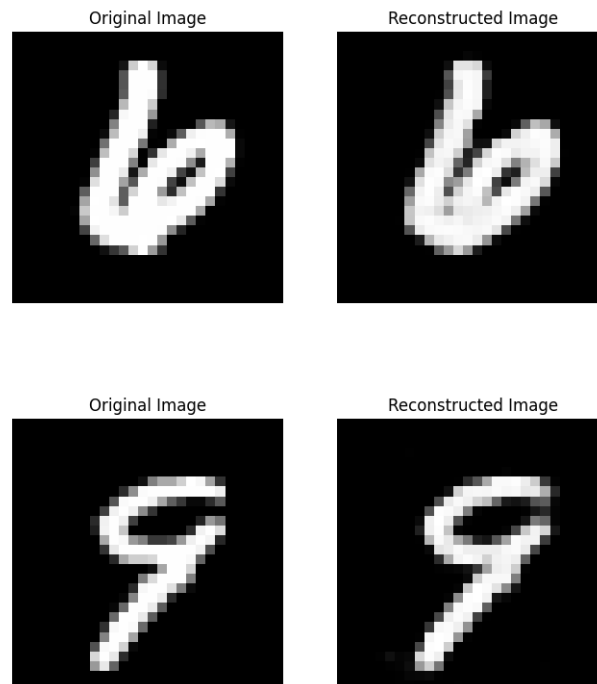
פיצ'רים ששוחזרו בהצלחה: המבנה הכללי של התמונה - התמונות המשוחזרות הגיעו לתוצאות כמעט מדויקות מבחינת הצורה של הספרה המקורית.

פיצ'רים שלא שוחזרו בהצלחה: הקצוות, הפרטים הקטנים - אפשר לראות בדוגמאות למטה, שהשוליים של הספרות מעט דהויים. בנוסף, זוויות חדות היו נקודת מכשול מבחינת שחזור.

זה הגיוני כי ה- GAN הזה אומן משלב התחלתי, ובשלב ההתחלתי D מתבסס על פריטים מאוד גסים כדי ליצור קלסיפיקציה, ולכן G מייצר תמונות שהפרטים הגסים שלהם יטעו את D . בשלב מתקדם יותר, תאורטית היינו יכולים להמשיך לאמן את G עם מפריד D' שהוא יותר חזק D וידע להבחין בפרטים ספציפיים יותר כדי לגרום ל- G ללמוד את הפרטים הקטנים בתמונה.

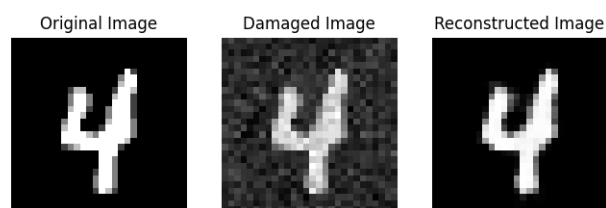
דוגמאות:

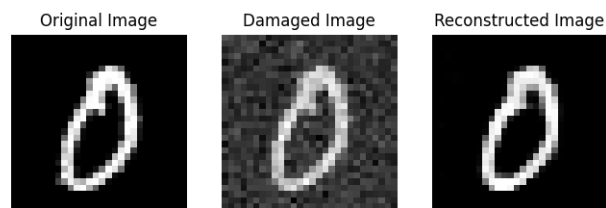
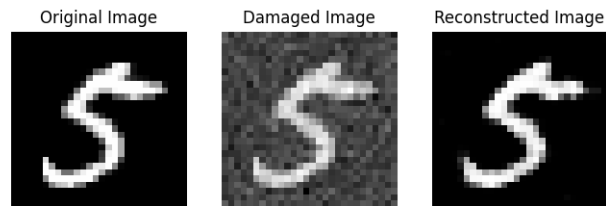




שאלה 3 Image Restoration:

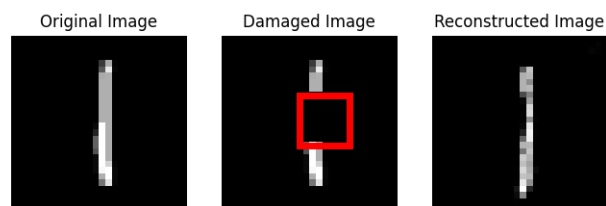
1. *Denoising*: בשביל משימה זו בחרתי את פונקציית לוס L_2 , ואלה התוצאות:

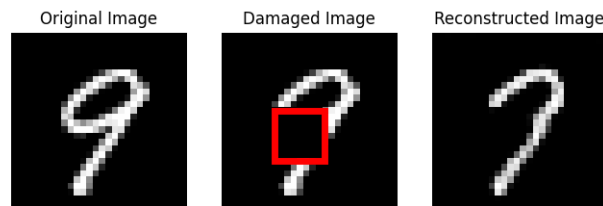
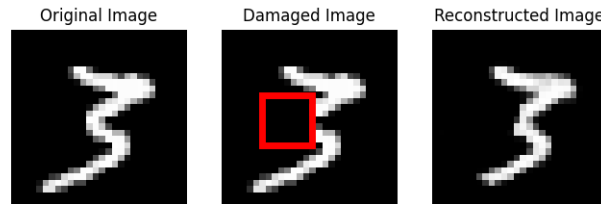




תוצאות מרשימות

2. *Inpainting*: למשימה זו בחרתי בL1. המשימה הזו הייתה לו הרבה יותר קשה....
תוצאות:





כמו שניתן לראות לפעמים הוא עבד טוב ולפעמים לא. הדוגמא של הספרה "1" היא לא יצאת דופן - בספרה הזו הוא הצליח יחסית לשחזר טוב, וזה הגיוני כי צורה מאוד גסה ללא הרבה פרטים קטנים. בשחזור של הספרה "9" יש עוד דוגמא לז פרטים, שאחד מהם גס והמודל משחזר, ואחד מהם עדין והמודל לא משחזר: הרגל הספרה 9 היא פרט גס, התבנית של קו ישר חוזרת על עצמה בהרבה ספרות ולכן המודל מצליח לשחזר את הרגל של הספרה. לעומת זאת, החלק המעוגל של ה-9 הוא פרט עדין, והמודל לא מצליח לשחזר אותו. ואכן, גם אם בן אדם היה מסתכל על התמונה הפגומה - זה די קשה להבין אם המשורר התכוון ל-7 או ל-9 או אפילו ל-1.

חלק תאורטי

שאלה 1

נסמן את התפלגות המקור ב- P_x , והתפלגות פשוטה שידועה לנו ב- P_z , כך ש- $x \sim P_x$, ונסמן ב- M את פונקציית החלפת המשתנה בין P_z ל- P_x , כך שמתקיים:

$$P_x(x) = \left| \frac{dM^{-1}(x)}{dx} \right| P_z(M^{-1}(x)), \quad M(z) = x$$

כש $\frac{dM^{-1}(x)}{dx}$ הוא היעקוביאן של M^{-1} , ו- $|\cdot|$ היא דטרמיננטה.

$GLOW$ משתמש ב- M בצורה דומה ל- GAN , אבל מניח עליו מספר הנחות, כדי שיוכל לחשב אותו ישירות. בהינתן דאטא, ותחת ההנחות הללו, אנחנו יכולים לאמוד את P_x על ידי חישוב ה- MLE בעזרת הנוסחה הנ"ל והדאטאסט. ההנחות שמתקיימות תחת $GLOW$ הן:

1. כל אחת מהשכבות ברשת היא הפיכה וניתן לחשב אותה בצורה אנליטית ומדויקת.
2. הדטרמיננטה של היעקוביאן של כל שכבה ניתן לחישוב בצורה יעילה (עד $O(n^3)$). מכאן נגזר תנאי שכל היעקוביאנים צריכים להיות מטריצות משולשיות (או פחות, מבחינת סיבוכיות).
3. המימדים בין השכבות אינם משתנים (אחרת, הפונקציות יהיו לאו דווקא הפיכות) נסמן את השכבות ברשת $h_1, h_2 \dots h_k$, כך שכל השכבות הפיכות ומתקיים:

$$M = h_k \circ h_{k-1} \circ \dots \circ h_1$$

נשים לב כעת, שמכך ש $\det(A^{-1}) = \det(A)^{-1}$ נובע ש:

$$P_x(x) = \left| \frac{dM(x)}{dx} \right|^{-1} P_z(M^{-1}(x))$$

בנוסף, מכיוון ש M היא הרכבה של פונקציות הפיכות אז גם M היא הפיכה ומתקיים $M^{-1}(x_i) = z_i$ עבור i סמפל בדאטאסט. אנחנו כאן מדברים על כל ההתפלגות, ולכן נסמן פשוט:

$$P_x(x) = \left| \frac{dM(x)}{dx} \right|^{-1} P_z$$

כדי לאמוד את P_x , נמקסם את הנראות שלה על פני הדאטא. ה- $\log - likelihood$ נתון על ידי:

$$\begin{aligned} \log(P_x(x)) &= \log \left(\left| \frac{dM(x)}{dx} \right|^{-1} P_z \right) \\ 1 &= \log \left(\left| \frac{dM(x)}{dx} \right|^{-1} \right) + \log(P_z) \\ 2 &= \log(P_z) - \log \left(\left| \frac{dM(x)}{dx} \right| \right) \\ 3 &= \log(P_z) - \log \left(\left| \frac{dh_k(z_{k-1}) \circ h_{k-1}(z_{k-2}) \circ \dots \circ h_1(x)}{dx} \right| \right) \\ 4 &= \log(P_z) - \log \left(\left| \frac{dh_k(z_{k-1})}{dh_{k-1}} \right| \cdot \left| \frac{dh_{k-1}(z_{k-2})}{dh_{k-2}} \right| \dots \left| \frac{dh_1(x)}{dx} \right| \right) \\ 5 &= \log(P_z) - \sum_{j=1}^k \log \left(\left| \frac{dh_j(z_{j-1})}{dh_{j-1}} \right| \right) \end{aligned}$$

כאשר 1 ו-2 מחוקי לוגריתמים, 3 הצבה של M כהרכבה של פונקציות, 4 מכללי דטרמיננטה, 5 מחוקי לוגריתמים. יש לשים לב שכל שכבה h_j פועלת על אינפוט מהשכבה הקודמת,

z_{j-1} ועל כן חישובי היעקוביאן בהתאם. כמו כן נגדיר את h_0 להיות פשוט $h_0 = x$ בשביל קונסיסטנטיות בסימונים.
אם כן, ה- MLE הוא:

$$MLE_{P_x} = \arg \max_{\{h_1, \dots, h_k\}} \left\{ \log(P_z) - \sum_{j=1}^k \log \left(\left| \frac{dh_j(z_{j-1})}{dh_{j-1}} \right| \right) \right\}$$

הרכיב הראשון ניתן לחישוב כי אנחנו יודעים את P_z , והחישוב של הרכיב השני כמובן נעשה ע"י שיטות מבוססות גרדיאנט.
המקרה שונה ב- GAN וב- GLO , שהם אינם ממדלים את ההתפלגות P_x אלא:

- GAN מנצל את **קיומו** של החלפת משתנה M שמחליף בין P_x ל- P_z , ורק משתמש באומדן שלו על מנת לייצר דגימות חדשות. בפרט, הוא לא מניח שום דבר לגבי M , ולא מחשב בעזרתו את P_x .
- GLO לא מנסה למדל את P_x , אלא מנסה לאפטם את $Range(G)$, בעזרת הוקטורים $latent\ space$, כלומר לגרום לתמונה של $Generator$ להיות הכי קרובה לדאטא סט, בעזרת מטריקות על z .

שאלה 2

בתחילת האימון של GAN , ל- D קל מאוד לזהות את התמונות המזויפות ש- G מייצר (זה בסך הכל לקבוע האם תמונה היא רעש) והוא קובע $Decision\ plane$ מאוד חד - כלומר עבור התמונות המזויפות הערך יהיה מאוד קרוב ל-0, ועבור אמיתיות מאוד קרוב ל-1. בנוסף, כל התמונות יהיו לו בבירור מזויפות כי G עוד לא ידע לייצר תמונות טובות. באשר ל- G , אנחנו מעדכנים אותו לפי הגרדיאנטים של D , ולכן בשלב הזה אם D ייתכנס - משמע, הגרדיאנטים שלו יהיו 0 - אז G לא יוכל ללמוד ונגיע לסיטואציה הסטורציה. להלן ההסבר האנליטי:

אנחנו מעוניינים לחשב את $\frac{\partial L(D)}{\partial G(z)}$, כלומר לחשב את הגרדיאנט של הלוס של D , ביחס ל- G . פונקציית הלוס היא כמו שלמדנו:

$$L(D) = \log(D(x)) - \log(1 - D(G(z)))$$

כש x תמונות אמיתיות ו- z וקטורים של רעש ב- $latent\ space$. בעקבות האמור לעיל לגבי מצב האימון ההתחלתי, נניח את ההנחה ש- $D(G(z)) = 1 - \varepsilon$, כלומר D יודע לסווג כמעט את כל התמונות המזויפות בהצלחה. נפתח את הביטוי:

$$\begin{aligned} \frac{\partial L(D)}{\partial G(z)} &= \frac{\partial \log(D(x)) - \partial \log(1 - D(G(z)))}{\partial G(z)} \\ &= - \frac{\partial \log(1 - D(G(z)))}{\partial G(z)} \\ &= - \frac{\partial \log(1 - (1 - \varepsilon))}{\partial G(z)} = - \frac{\partial \log(\varepsilon)}{\partial G(z)} \\ &= - \frac{1}{\varepsilon} \cdot \frac{\partial \varepsilon}{\partial G(z)} \\ &= - \frac{1}{\varepsilon} \cdot 0 = 0 \end{aligned}$$

כאשר 1 נובע מכך ש $D(x)$ לא תלוי ב G , 2 נובע מההנחה למעלה, 3 נגזרת של \log , 4 נובע מכך ש ε הוא קבוע ביחס ל $G(z)$.
 כלומר, אם D התכנס במצב ההתחלתי - ל G לא תהיה השפעה על הגרדיאנטים שכבר התאפסו, והלמידה תיעצר - סטורציה.

שאלה 3

בשאלה הזו אנחנו עוסקים במצב של *Mode Collapse*.

1. הבעיה הפנימית $\{V(D, G)\}$ מאפסמת את G ביחס ל D קבוע כלשהו. למעשה, היא מוצאת את המצב של G שגורם לתחזיות של D להיות מינימליות. הסבר: ניזכר בהגדרה של V כמו שהגדרנו בכיתה:

$$V(D, G) = \mathbb{E}_{p_{data}} [\log(D(x))] - \mathbb{E}_{p_z} [\log(1 - D(G(z)))]$$

אם נתייחס למצב בו D קבוע, אז ישנו ערך קבוע של G עבורו מתקבל $D(G(z))$ מינימלי (ובך כפועל יוצא $\mathbb{E}_{p_z} [\log(1 - D(G(z)))]$ מקסימלי ובכך לבסוף V מינימלי כמו שהבעיה הפנימית דורשת) והמשמעות היא שהדבר נכון **לכל** z , כלומר עבור כל אינפוט הוא יוציא את אותה התמונה שעבורה D יביא ערך מקסימלי. זוהי בדיוק הבעיה של *Mode Collapse*.

2. הבעיה החיצונית $\{Max_D \{Min_G \{V(D, G)\}\}$ מאפסמת את D שיביא ערך מקסימלי, בהינתן G שכבר התכנס, או לצורך העניין - G מושלם, שמייצר תמונות שנראות לגמרי אמיתיות. זו משימה קשה עבור רשת מוגבלת מבחינת פרמטרים, כי הבחנה בין תמונה אמיתית לתמונה מזויפת מצויינת - תלויה בפרטים מאוד מאוד קטנים, ולצורך כך צריך D חזק ובפרט גדול. בנוסף, ככל שהרזולוציה של התמונה גדולה יותר כך צריך D עוד יותר גדול שיוכל לתפוס את כל הקשרים והפרטים בתמונה.

3. אני רואה בסנריו הזה 2 סכנות - כאמור, הראשונה היא *Mode Collapse* כמו שהוסבר למעלה.
 שנית, במקרה ש D לא מצליח להפריד טוב בין תמונות אמיתיות לתמונות מזויפות, הגרדיאנטים שהוא ייצר יהיו בכיוון לא נכון עבור האימון של G , ובעצם הדבר יוריד את הביצועים הכלליים.