

# סיכום רשתות נוירונים לתמונות

11 באוגוסט 2023

## 1 רשתות קונבולוציה - CNN:

אנחנו משתמשים ברשתות קונבולוציה לתמונות משום שהן עובדות טוב עם שתי ההנחות הבאות:

1. כל נוירון לא מושפע מכל התמונה, אלא רק חלק קטן שנקרא - receptive field.

2. בתמונות אנחנו משתמשים בסטציונריות מרחבית, כל פריט יכול להופיע במקומות שונים ובצורות שונים בתמונה, ואנחנו רוצים שהמודל ינתח את התמונה על כולה, ולא לפי מיקום ספציפי - translation variance.

נשים לב כי בתמונות נרצה שכל הנוירונים באותה שכבה יגיבו באותו האופן לאותם גירויים (כי האובייקט יכול להופיע במקומות שונים התמונה), לכן הם יהיו עם אותם משקלים weight shering.

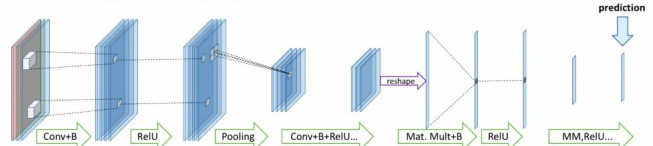
**רשת קונבולוציה:** רשת שתכיל שכבות כך שכך שכבה מקבלת את אותו המשקל, כל שכבה תפעיל פולטר שנלמד על הקלט. היא תעשה את זה על חלון בגודל מסויים שאנחנו נקבע מה גודלו.

**למעשה:** הרשת מזהה פרמטרים קטנים בתמונה ועושה להם לבסוף אינטגרציה, ובודקת האם היא מכילה אובייקט מסויים.

### 1.1 ארכיטקטורה של Alex Net:

#### Prototypical Classification CNN Architecture

- Consists of repeated layers of:
  - Convolution** – linear translation-invariant (LTI) operators with a restricted spatial support (compact filters)
  - Bias** – addition of activation offset (threshold)
  - Pointwise Non-linearity** – modeling the neural activation
  - Pooling/Striding** – downscaling the spatial signal resolution
  - Fully-connected (MLP)** – at coarse resolution, all neurons are connected



### 1.1.1 שכבת קונבולוציה:

אוסף של טרנספורמציות לינאריות - לינארי ואינווריאנטי לזמן או למרחב.

$$(Lx_{i+k})_j = (Lx_i)_{j+k}$$

הם לינארים ולכן ניתן לכתוב אותה כמטריצה, כך שהמטריצות יקיימו תכונות סויימות. מטריצה סימטרית עם אלכסונים שווים, ויש בה מלא אפסים.

הקשר בין קונבולוציה להתאמת דפוס - pattern matching:

$$\sum_x (I(y+x) - f(x))^2 = \|I(y+x)\|^2 + \|f(x)\|^2 - 2 \sum_x (I(y+x)f(x))$$

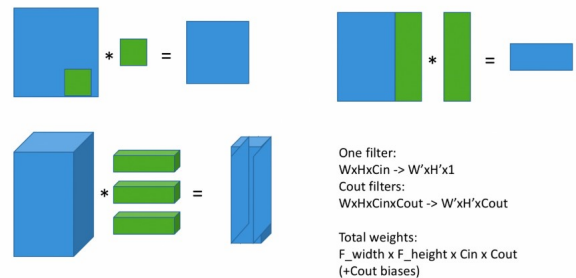
$$\sum_x (I(y+x)f(x)) = \sum_x (I(y-x)f(-x)) = (I * \bar{f})(y)$$

למעשה, הרשת מחפשת דפוסים בדאטה, אך לא בהכרח דפוס שאנחנו חשבנו שהוא חשוב.

- אנחנו מורידים את הרזולוציות המרחביות, ולכן כדי לפצות על זה אנחנו נעלה את מספר הערוצים (פיצ'רים). כך אנחנו יותר אבסטרקטים.
- ה receptive field של הנוירונים בשכבות האחרונות גדול יותר, כי הוא מכיל את כל התמונה כמעט.
- בתמונות טבעיות, יש קורולציות מרחביות במרחקים כטנים, לאן אין טעם להשתמש בפילטר גדול, כי ככל הנראה לא יהיה קשר בין חלקים רחוקים בתמונה.
- הקונבולוציה היא על הצירים  $x, y$ . אך על ציר ה  $z$  אנחנו משתמשים ב"פ", ולכן אנחנו מקבלים סקלר. כלומר - כל פילטר מופעל על  $x, y$ , וציר ה  $z$  מצטמצם ל 1. כל פילטר מחזיר לנו ערוץ בודד לשכבה הבאה. מספר הפילטרים בשכבה ה  $i$ , יהיו מימד ה  $z$  של הקלט בשכבה ה  $i + 1$ .
- נשים לב כי התמונה שנקבל לאחר פילטר תהיה קטנה יותר בגלל קצוות התמונה.

תוצאות של הפעלת פילטרים שונים על תמונות שונות:

#### Convolutional Layer



### 1.1.2 שכבה לא לינארית והביאס:

- נרצה להכניס למודל פונקציה לא לינארית כדוגמת: סייגמואיד, או  $ReLU$ .

$$Sigmoid := \frac{1}{1 + e^{-x}}$$
$$ReLU := \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$$

- בחירה של פונקציה ספציפית לא כלכך תשנה, העיקר שתהיה פונקציה שמממשת אי לינאריות.
- לבסוף נוסיף  $bias$  שמזיז לנו את הסף של הגרף. ביאס חיובי יזיז את הגרף שמאלה, וביאס שלילי יזיז את הגרף ימינה.

### 1.1.3 שכבת Pooling:

- שכבה בה אנחנו נוריד את מימד התמונה.
- יש לנו מודל שנקרא  $stribe$  שבו אנחנו נדגום כי פיקסל שני אם  $stribe = 2$  וכן הלאה.
- ניתן להשתמש ב  $avrage pooling$  או  $Max pooling$ .
- אין טעם ללמוד  $Average$  כי הרשת הייתה יכולה ללמוד את זה בעצמה אם זה היה חשוב.
- $Max$  יהיה טוב יותר כי הוא לא לינארי.

## 1.2 השכבה האחרונה ומדידת ביצועי הרשת:

ברשת קלאסיפיקציה יש לנו מספר קלאסים שאנחנו רוצים שהרשת תמפנ אליהם את הקלט. כלומר בהינתן קלט, נרצה שהרשת תגיד לנו מה הקלאד שהקלט שייך אליו.

### 1.2.1 Softmax:

כדי לבצע זאת, נוסף לאחר שכבת ה  $FC$ , שכבת  $softmax$ . הרשת פולטת לנו ווקטור לאחר שכבת ה  $FC$  (משום שכפלנו מטריצה בווקטור), מימד הווקטור יהיה מספר הקלאסים. אנו נרצה להמיר את הווקטור שהרשת פולטת לווקטור הסתברות, שמנבא את ההסתברות שהפלט שייך לכל אחר מהקלאסים. נעשה זאת בעזרת פונקציית  $softmax$ :

$$p_i = \frac{e^{q_i}}{\sum_i e^{q_i}}$$

לאחר מכן נמדוד את ביצועי הרשת בעזרת  $Loss$ . כמה רחוק הווקטור שיצא משכבת ה  $softmax$  לעומת הווקטור המקורי  $lable$  -

## 1.2.2 פונקציות Loss:

### 1. Cross Entropy Loss - לרשת קלאסיפיקציה:

נמדוד את ביצועי הרשת באופן הבא:

כאשר  $p$  הוא ווקטור הפלט ו  $y$  הוא ה  $label$ .

$$H(p, y) = H(p) + D_{KL}(p||y), \quad H(p, q) = - \sum y_i \log(p_i) \quad D_{KL}(p || y) \equiv p_i \log \left( \frac{y_i}{p_i} \right)$$

למעשה, אנחנו רוצים פונקציה גזירה כדי שנוכל להביא אותה למקסימום, או להביא את ההפסד למינימום.

אם אנחנו מתמודדים על שני קלאסים ב"ת, כלומר שקיום של קלאס אחד מעיד בהכרח שהשני לא מתקיים. נוכל למדוד באופן הבא:

$$2 \text{ classes} : L(p, y) = -y_i \log(p_i) - (1 - y_i) \log(1 - p_i)$$

$$\text{more classes} : L(p, y) = - \sum_i y_i \log(p_i)$$

כאשר את מיפוי הערך לטווח  $[0, 1]$  נוכל לעשות בעזרת פונקציה לוגיסטית:

$$\text{logistic function } p = \frac{1}{1 + e^{-q}}$$

### 2. Regression Loss - לרשת רגרסיה:

בבעיות רגרסיה אנחנו רוצים שהרשת תחזה לנו מספר בטווח (למשל גיל או מחיר), כדי למדוד את ביצועי הרשת נוכל להשתמש בפונקציית ה  $Loss$  הבאה:

$$\sum_n (N_\theta(x^n) - y^n)^2$$

## 1.3 אימון הרשת - למידה:

אנו רוצים ללמוד את משקולות הרשת, כדי שהיא תחזה לנו את הפלט הנכון. למעשה אנחנו רוצים למזער את ה  $Loss$ , נעשה זאת ע"י פונקציות גזירות ומציאת המינימום ע"י הנגזרת.

$$\min_{\theta} loss = \frac{\min_{\theta} - \sum_n \sum_i y_i^n \log p_i^n}{\sum_n (N_\theta(x^n) - y^n)^2}$$

## 1.4 האלגוריתם: Gradient Descent Optimization

יש לנו פונקציה, ואנחנו מנסים למצוא את המינימום \ מקסימום שלה. נעשה זאת ע"י הליכה בכיוון הגרדיאנט (השינוי הגדול ביותר בנגזרת בכל אחד מהכיוונים  $x, y$ ) לחלופין עבור בעיית מינימום נלך בכיוון המנוגד לגרדיאנט. נבחר את גודל הצעד, וכך בכל שלב נלך בכיוון הגרדיאנט בגודל הצעד שבחרנו, עד שניתכנס למינ' \ מקס' מקומי (כשהגרדיאנט יתאפס).

פורמלית:

### Gradient Descent Optimization

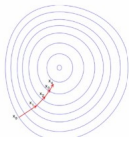
The gradient points to the direction of the maximum change:

$$\max_{\|z\|^2=\varepsilon} f(x+z) - f(x) = \max_{\|z\|^2=\varepsilon} f(x) + z \cdot \nabla f(x) - f(x) + O(\varepsilon^2) = \max_{\|z\|^2=\varepsilon} z \cdot \nabla f(x) + O(\varepsilon^2)$$

$$\text{Lagrange multipliers} \quad \max_{\|v\|^2=1} v \cdot w \Rightarrow \nabla_v (v \cdot w) = \lambda \nabla (v^T \cdot v) \Rightarrow w = \lambda v \Rightarrow v = w / \|w\|$$

Minimization of general (mostly) differentiable, is often done via a gradient descent process

$$\theta^{(k+1)} = \theta^{(k)} - \nabla_{\theta} \text{loss}(N_{\theta^{(k)}}, x, y)$$

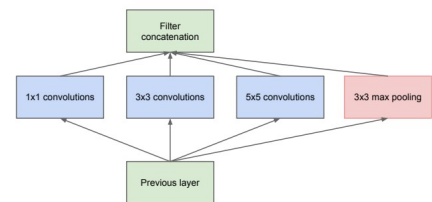


## 2 ארכיטקטורת רשתות:

- חוקרים מצאו כי עדיף להחליף שכבת קונבולוציה עם קרנל גדול, במספר שכבות קונבולוציה עם קרנלים קטנים יותר. כך אנו מקטינים את מספר הפרמטרים, וגם יש לנו יותר שכבות אי לינאריות שנותנות אקספרסיביות לרשת.
- Inception Module ברשת *GoogLeNet*: חוקרים של גוגל החליטו לפצל את הרשת לרוחב, כך שבשלב מסויים יופעלו עליה כמה קרנלים בגדלים שונים, בנפרד. לאחר מכן נחבר את כולם, כך שהרשת תבחר את הקרנל הטוב ביותר.

### Inception Module

- Tough decisions: convolution or pooling? what size of filters to use?
- Idea: do multiple options in parallel!



- קונבולוציה עם קרנל  $1 \times 1$ : למעשה היא מעבירה ווקטור לסקלר. אם נפעיל קונבולוציה עם קרנל כזה על מטריצה  $n \times n \times m$  אזי כל קאורדינטה  $(i, j)$  תעבור להיות סקלר. כך למעשה ניתן להוריד או להעלות מימד, כתלות במספר הפילטרים בקונבולוציה  $1 \times 1$ .
- נרמול האינפוט - Batch Normalization: בשלב מסויים גילו כי יעיל לנרמל את האינפוט של כל שכבה באופן הבא:

$$\text{BN}(x) = \gamma \left( \frac{x - \mu(x)}{\sigma(x)} \right) + \beta$$

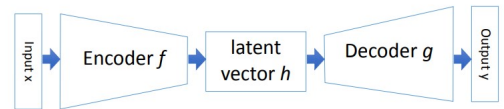
כך שהתוחלת תהיה 0, ותהיה שונות מסויימת.  
 באוסן כזה אין הבדלים בין התמונות, אלא רק השוני שאחנו רוצים שהרשת תלמד.  
 $\gamma, \beta$  הם פרמטרים נלמדים, ו  $\sigma, \mu$  הם ווקטורים.

- הרשת *ResNet*: רשת זאת של מייקרוסופט באה לפתור את בעיה אימון הרשתות העמוקות. עד אז, רשתות עמוקות לא חזו טוב, והרשת הזאת עושה זאת ע"י זה שהשכבות לומדות משהו מסויים, לאחר מכן מוסיפים להן את הקלט המקורי. כלומר, בשלבים מבויימים מכניסים לרשת את האינפוט לפלט של השכבה האחרונה. כך ערכי הגרדיאנט לא יהיו קטנים מידי ונוכל לחשב אותם ביותר יעילות.

### 3 *Autoencoders*

רשת שמבוססת על *encoder – decoder*.  
 ניתן להשתמש בה להורדת מימד לא לינארית, או להסקת אינפורמציה על תמונות משום שהיא שומרת מטא דאטה שיכול לשמש ללמידה. בנוסף הרשת משמשת למודלים גנרטיביים. רשת זאת היא רשת ללמידה **בלתי מפוקחת**.

**מבנה הרשת:** הרשת מקבלת אינפוט ומקודדת אותו לוקטור במימד מסוים. בנוסף יש רשת *decoder* שמקבלת ווקטור מקודד ומוציאה את הפלט המקורי (מטרתה לשחזר את האינפוט).



כך, אנחנו יכולים להוריד מימד, כי רשת לומדת את הפרטים החשובים ביותר.

**חישוב ה *Loss*:** נמדוד עם נורמה 2 כך:

$$L(x, y) = \sum \|y - g(f(x))\|^2 = \sum \|x - g(f(x))\|^2$$

**טענה:** אם ה *encoder* - מטריצה לינארית, וה *decoder* הוא אותה מטריצה בטרנזספוז - אז מתקיים כי הרשת מחשבת את המטריצה של *PCA* (אלגוריתם לינארי להורדת מימד).

**כדי למצוא את ההתפלגות נשתמש בלוג לייקליהוד:** נוכל להחליף את זה בנורמה 2

$$L(x, y) = -\log P(y | h(x)) = -\log P(x | h(x))$$

**שימוש ב *cross – entropy*:**

$$L(x, y) = -\sum_i x_i \log(y_i) + (1 - x_i) \log(1 - y_i)$$

**Sparse Autoencoders:** אנחנו נרצה לגרום לווקטור הפלט של ה *encoder* להיות דליל. נעשה זאת כך ע"י הוספת גורם נוסף ללוס -

$$L(\mathbf{x}, g(f(\mathbf{x}))) + \Omega(\mathbf{h})$$

ע"י הוספת עונש  $\Omega(\mathbf{h})$  ששווה ל:

$$\Omega(\mathbf{h}) = \lambda \sum_i |h_i|$$

או העונש הבא, שיציב יותר לרעש ולשינויים קטנים באינפוט:

$$\Omega(\mathbf{h}, \mathbf{x}) = \lambda \sum_i \|\nabla_{\mathbf{x}} h_i\|^2$$

**Denoising Autoencoders:** נקח את האינפוט המקורי  $x$  ונוסיף לו רעש באופן מלאכותי ונגדיר את האינפוט המורעש להיות  $\tilde{x}$ . לאחר מכך נדרוש מהרשת שלנו להוציא פלט שדומה ל  $x$ . כלומר הרשת תלמד לנקות רעשים מתמונה.

$$L(\mathbf{x}, g(f(\tilde{\mathbf{x}})))$$

## 4 מודלים גנרטיבים - *GANs* - Adversarial Generative Models

**המשימה:** המודל צריך ללמוד את כל ההתפלגויות שיש לנו, בניגוד לרשתות קלסיפיקציה. מודלים אלו יכולות לייצר דגימות מתוך ההתפלגות - לייצר תמונות חדשות.

**משימה נוספת:** מודל זה יכול לדגום תמונה בהינתן משט, או תמונה מטושטשת.

**כיצד ניצור תמונה:** נדגום ווקטור  $z$  מההתפלגות הרב מימדית שלנו, לאחר מכן נכניס את הווקטור  $z$  לפונקציה  $M_\theta$  שתמפה אותו, כך ניצור מ"מ חדש.

$$M_\theta(z) = x$$

כאשר  $x \sim P(x)$ , אך אנחנו לא יודעים מהי ההתפלגות  $P(x)$ .

למעשה  $M_\theta$  תהיה רשת נוירונים מרובת פרמטרים, ואנו נלמד את הפרמטרים שלה.

#### Sampling by Change of Variable

Map easy-to-sample dist. (e.g., Gaussian, Uniform) into a target dist.

$$M_\theta(z) = x, \quad z \sim N(0, I), \quad x \sim P(x)$$

change-of-variable      easy-to-sample dist.      target distribution

$$P(x) = \left| \frac{\nabla M^{-1}(x)}{\nabla z} \right| P_z(M^{-1}(x))$$

better suited for sampling than opt.  $P(x)$       density expansion/contraction term

## 4.1 כיצד נאמן את $M$ :

יש לנו שתי רשתות  $G$  – generator שהיא למעשה תהיה  $M$ , ורשת  $D$  – discriminator, שתי הרשתות ילמדו אחת את השנייה.

- $G$ : תקבל  $z$ -ים מההתפלגות הפשוטה שאנחנו מכירים, לאחר מכן  $G$  תייצר לנו תמונה.
- $D$ : תקבל שני פרמטרים - את הפלט של  $G$ , ותמונה מהדאסט סט. היא תצטרך להגיד איזה מהתמונות אמיתית (1) ואיזו מג'נרטת ע"י  $G$  (0).
- המטרה של  $D$  היא למצוא את התמונה המזוייפת, והמטרה של  $G$  היא להטעות את  $D$ .

האימון של שתי הרשתות נעשה באמצעות פונקציית הלוס הבאה:

$$\min_G \max_D V(D, G) = \min_G \max_D E_{p_{\text{data}}} [\log D(x)] + E_{p_z} [\log(1 - D(G(z)))]$$

קודם כל, נרצה ש  $D$  תתכנס לאנשהו - ניתן לו מלא איטרציות, ולאחר מכן נתחיל להריץ את  $G$  עם מעט איטרציות. **מתי נסיים:** כשהרשת  $G$  תצליח להערים על  $D$ , אח"כ נזרוק את  $D$  ונשתמש ב  $G$  שתייצר לנו תמונות.

## 4.2 בעיית Mode Collaps:

- **בעיית סטורציה:** אם  $D$  מזהה טוב מידיי אילו תמונות הן המזוייפות, אזי  $G$  לא יוכל לגזור את הפונקציה ש  $D$  מחזיר לו - הוא יגיע לסטורציה (הגרדיאנטים יתאפסו). ולכן לא משנה מה הוא יעשה הוא לא יצליח לרמות את  $D$ .
- הפתרון:** נעזור ל  $G$  - במקום לתת ל  $D$  להגיע להתכנסות וכך הוא יצבור יתרון על  $G$ , נריץ כל אחר מהם לתהליך של גרדיאנט דסנט אחד.
- הפתרון הקודם יכול לגרום לבעיה אחרת: אם ננסה למזער את  $G$  קודם, במקום למקסם את  $D$  תחילה. כך -

$$\min_G \max_D V(D, G) \Leftrightarrow \max_D \min_G V(D, G)$$

אזי  $G$  יבחר תמונה אחת ש  $D$  חושב שהיא אמיתית, וישלח לו אותה בכל פעם כדי "לנצח".



**לכן:** נצטרך לאמן בזהירות ולבחור את ההיפר פרמטרים כך שלא ניכנס לבעיית Mode Collaps

**בעיות נוספות של GAN:** הם לא יודעים לייצר תמונות מהדאטה סט, כלומר הם לא יודעים לייצר את כל התמונות מההתפלגות.

### 4.3 Conditional GAN - דגימת התפלגויות מותנות:

שיטה זו נכונה באופן כללי, אך נראה אותה בהקשר של GAN. אנחנו דוגמים וקטור מתוך התפלגות  $p(x)$ , אך למעשה אנחנו יכולים להתנות אותה במאורע  $y$  כך  $p(x|y)$ . **לדוגמה:** שהגנרטור שלנו ייצר ספרה מסויימת.

#### הביצוע:

נכניס ל  $G, D$  את ה  $class$  כך שהם יידעו לאן לכוון.

- נצטרך להכניס לדסקרימינטור את התיוג של התמונה בכל פעם, כך הוא יראה את הקורלציה בין התמונה לתיוג.
- לעומת זאת הגנרטור יצטרך לייצר את התמונה באופן אקראי, וגם לדעת מה אנחנו רוצים שהוא ייצר.

#### הארכיטקטורה:

אנחנו צריכים לשנות את ארכיטקטורת הרשת, כי אנחנו רוצים להזין את ה  $class$  גם כן לתוך הרשת. ניתן לעשות זאת ע"י הכנסת כופל של ווקטור  $bias$  שנלמד.

#### המודל:

בהינתן מונה  $I_i$  של  $class$  שמתוייג  $C_i$ :

$$D(I_i, c_i) - > 1$$

$$D(G(z, c_i), c_i) - > 0$$

#### עבור העלאת רזולוציה של תמונות - SRGAN:

ניתן לאמן באופן הבא:

$$D(H_i) - > 1 \quad D(H_i) - > 1$$

$$D(G(z, L_i)) - > 0 \quad D(G(L_i)) - > 0$$

כאשר  $L_i = \text{down}(H_i)$  ו  $z$  הוא ווקטור. ניתן לאמן את הרשתות גם ללא הווקטור  $z$ .

דרך נוספת לייצר התפלגויות מותנות בעזרת GAN, עם ארכיטקטורה טובה יותר.

**דאטה סט:** יש לנו צמדים של תמונות ואנחנו רוצים ללמוד את ההתאמה בין הצמדים.

**לדוגמה:** איור של תמונה, ותמונה. ואנחנו רוצים שהרשת תקבל איור ותייצר ממנו תמונה.  
**נלמד כך:** כאשר  $I, J$  תמונות (איור ותמונה).

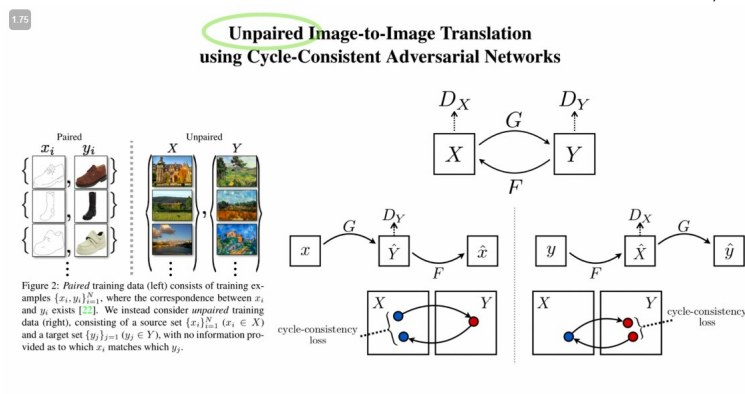
$$D(I_i, J_i) - > 1$$

$$D(I_i, G(I_i)) - > 0$$

$G$  מקבל את האיור ומייצר ממנו תמונה.

**כיצד נייצר זוגות של תמונות:**

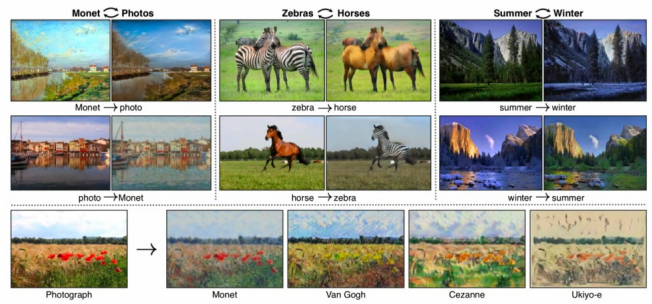
ניתן להשתמש בשני רשתות GAN כדי לייצר את הזוגות באופן הבא:



כך, נייצר מידע משני קלאסים -  $X, Y$ . נקח תמונות מ  $X$  נכניס ל  $G$  והוא יעביר אותנו ל  $Y$ . באופן האופן נמפה תמונות מ  $Y$  ל  $X$ .

נעשה זאת כך:  $G$  ימפה תמונה מ  $X$  ל  $Y$ , לאחר מכן נבקש ממנו למפות את התמונה הזאת חזרה ל  $X$ . באופן כזה נקבל את אותה התמונה (תיאורית). נעזה באותו האופן מ  $Y$  דרך  $X$  חזרה ל  $Y$ , ונרצה ששני המעגלים ימוזערו. כך יהיו לנו זוגות של תמונות.

**התוצאות נראות כך:**



## 5: Non-Adversarial Generative Models

שיטות שהתפתחו כאסטרטגיה חלופית לרשתות  $GAN$ . באותו האופן של  $GAN$  נשתמש בחילוף משתנה - אנחנו רוצים ללמוד התפלגות, והרשת שלנו צריכה ללמוד את ההתפלגות שנשלטת ע"י משתנה נלמד.

### 5.1: GLOW – Generative Flow

היוצרים בחרו רשת שהארכיטקטורה שלה תייצר לנו מטריצה משולשית הפיכה שקל לחשב את הדטרמיננטה שלה. כדי שנוכל לפתור את הבעיה הבאה:

#### GLOW: Generative Flow

Reminder:

$$M_\theta(z) = x, \quad z \sim N(0, I), \quad x \sim P(x)$$

change-of-variable      easy-to-sample dist.      target distribution

$$P(x) = \left| \frac{\nabla M^{-1}(x)}{\nabla z} \right| P_z(M^{-1}(x))$$

density expansion/contraction term

Requirements from  $h_k$ :

- invertible
- known inverse
- efficient Jacobian determinant
- maintain same dim.

In case of a composition:  
(as the case of multi-layered net.)

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

$$\log P_\theta(x) = \log p(z) + \log |\det(dz/dx)|$$

$$= \log p(z) + \sum_k \log |\det(dh_{k-1}/dh_k)|$$

simple dist.      can this be made simple?

triangular Jacobians:

$$\log |\det(dh_{k-1}/dh_k)| = \sum_j (\log |(dh_{k-1}/dh_k)_{jj}|)$$

Raanan Fattal

67103 Introduction to Neural Networks

24

### 5.2: GLO - Generative Latent Optimization

שיטה נוספת, שאינה אדברסריאלית.

**הרעיון:** נאפטם את  $G$  כך:

$$\min_{\theta_G, \{z_i\}} \sum_{i \in B} \|G(z_i) - I_i\|$$

כדי שיוכל לפרוש את ההתפלגות שלנו, אך היא מאמנת את פרמטרי הרשת ואת הפרמטרים של וקטורי המקורות.

נאתחל את ה  $z$ -ים להיות דגימות של גאוסיאן סטנדרטי. כך ההתפלגות של  $z$  לא תזוז הרבה, ולכן נוכל לדגום מהגאוסיאן הרגיל.

## GLO: Generative Latent Optimization

[Optimizing the Latent Space of Generative Networks]

Encoder-less autoencoder:

- Replace the encoder by free (learnable) codes  $z_i$  (one for each example)
- Equiv. to infinite capacity encoder (up to initialization)

Solve:  $\min_{\theta_G, \{z_i\}} \sum_{i \in B} \|G(z_i) - I_i\|$

- Unlike in AE, latent vectors can be initialized as desired

Using as a Generator:

- initialize  $z_i$  to simple known distribution
- Minimize changes in distribution using low LR for  $z_i$
- assume it didn't, sample  $z$  from initial dist.

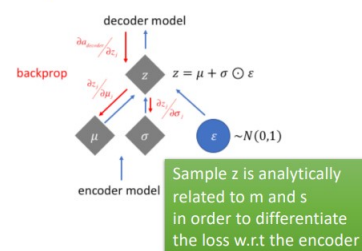
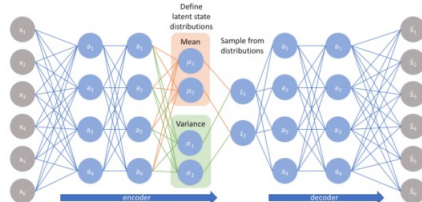
## 5.3 Variational Auto-Encoders (VAEs)

השיטה הזאת גם לא ככה טובה כמו הקודמות.

**מבנה הרשת:** יש לה מבנה של אוטו-אינקודר. אנחנו מכניסים תמונות ומקבלים קידוד. אל הווקטור שאנחנו מקבלים ב  $latent space$  אנחנו נתייחס כאל ווקטור התפלגות שנדגם מבתפלגות מסויימת, האנקודת מייצר את התוחלת והשונות. כלומר התמונה הופכת להתפלגות.

לאחר מכן כנניס את התמונה ל  $decoder$  שייצר לנו תמונה מהווקטור.

## Variational Auto-Encoders (VAEs)



Wants latent space to be a unit Gaussian

Every input point (image) is mapped to a Gaussian dist. in latent space

- Encoder outputs the mean and variance of the Gaussian
- A sample is drawn from this Gaussian
- Decoder reconstructs the input

**שלב האימון:** כל תמונה מתמפה לגאוסיאן ספציפי לתמונה. דוגמים מהגאוסיאן לפי הנוסחה שמופיעה בשקף  $z = \mu + \sigma \cdot \epsilon$  מכניסים ל  $decoder$  ומגדירים  $Loss$  שישחזר את התמונה.

**הדגימה:** יש בשיטה שני לוסיים, אחד של האוטו-אינקודר, ואחד נוסף. אנו רוצים שכל אחד מהגאוסיאנים יהיו סטנדרטי (תוחלת 0 ושונות 1) משום שכל תמונה מתמפה לגאוסיאן אחר, אך בסופו של דבר כדי לייצר תמונות אנחנו רוצים לדגום מגאוסיאן אחד, ולשלוח ל  $decoder$ . לכן נרצה לוס נוסף שימדוד לנו את הגאוסיאן הכללי.

**פונקציית הלוס:** נשתמש בלוס הבא

$$D_{KL}(G_1, G_2) = \int (\log(G_1(x)) - \log(G_2(x))) G_1(x) dx =$$

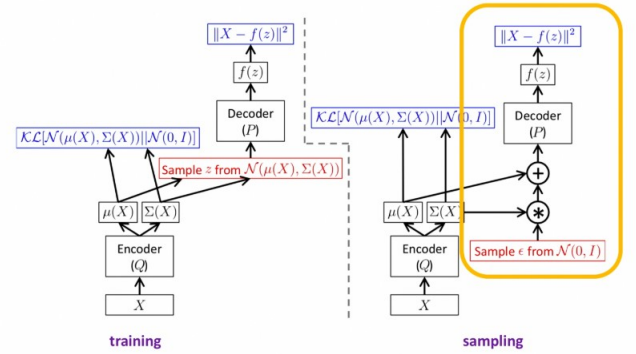
$$\int \left( \sum_k (x^k - \mu_1^k)^2 / 2\sigma_1^k - \sum_k (x^k - \mu_2^k)^2 / 2\sigma_2^k \right) (\sqrt{2\pi^d}) e^{-\sum_k (x^k - \mu_1^k)^2 / 2\sigma_1^k} dx$$

הלוס הסופי:

$$\sum_i \left\| \text{reconstruction} \left( \left\| (I_i) \right\| + E_{D_7} (I_i)^2 - E_\sigma (I_i)^2 - 1 - 2 \log (E_\sigma (I_i)) \right) - I_i \right\|$$

סיכום:

## VAEs Summary



## :Wasserstein Auto-Encoders 5.4

שיטה זאת בעיה לפתור את הבעיות בשיטה הקודמת  
 בסיטה הקודמת דרשנו כי כל גאוסיאן יהיה סטנדרטי, כדי שהגאוסיאן הכללי יהיה סטנדרטי. אך זה גורם לכך שתמונות שונות ימופו לאותה תמונה משום שה *latent vector* לקח אותם לגאוסיאנים שונים.  
 הפתרון הוא להוסיף לוס נוסף.

פונקציית הלוס:

$$\text{MMD}[\mathcal{F}, p, q] := \sup_{f \in \mathcal{F}} (\mathbf{E}_{x \sim p}[f(x)] - \mathbf{E}_{y \sim q}[f(y)]) , .$$

$$\text{MMD}[\mathcal{F}, X, Y] := \sup_{f \in \mathcal{F}} \left( \frac{1}{m} \sum_{i=1}^m f(x_i) - \frac{1}{n} \sum_{i=1}^n f(y_i) \right)$$

נקח הרבה תמונות, ונחשב את התוחלת של פונקציה ביחס לדגימות הללו.  
 לאחר שהאנקודים ייצרו לנו מהתמונות *latent vector*, נקח אותם. נחשב את ממוצע שלהם על איזושהי פונקציה  $f$ , ונשווה לדגימות שנדגמו מהתפלגות נורמלית סטנדרטית. כך נוכל לבחור את הפונקציה  $f$  המתאימה ביותר כך שהאנקוד ימצא לנו את ההתפלגות הסטנדרטית של הדגימות.

## 6 :Score Matching / Denoising Diffusion Probabilistic Models

**בסיס השיטה - SDE:** בשונה משיטות קודמות, כאן לא נשתמש בשינוי משתנה.

$$\frac{\partial x}{\partial t} = \nabla E(x)/2T + \eta$$

$$\frac{\partial}{\partial t} x \approx \frac{x^{n+1} - x^n}{\Delta t}, \quad x^{n+1} = x^n + \Delta t \frac{\nabla E(x)}{2T} + \sqrt{\Delta t} \eta, \quad \eta \sim N(0, 1)$$

הרעיון הוא להסתכל על משוואה דפרנציאלית שיש לה חלק דטרמיניסטי שסוכם את סך כל ההשפעות והשינויים עד הזמן  $t$  הנוכחי, ועוד הוספת רעש. באופן כזה אנו יכולים לצאת מנקודה מסויימת, ולהגיע להתפלגות. את המשתנה  $t$  אנחנו מגדירים, והוא מסמן את גודל הצעד. המשתנה  $T$  הוא מעין פרמטר רגולריזציה שמגדיר מ משקל הרעש לעומת משקל פונקציית האנרגיה  $E$ . **בסוף נתכנס ל:**

$$P(x) \propto e^{-E(x)/T}$$

אנו רוצים ללמוד את פונקציית האנרגיה  $E$ .

**השיטה:** נתחיל מ  $x^0$ , וננסה להגיע לנקודת שוויון. כשנגיע לשוויון זאת אומרת שהתייצבנו על התפלגות שחיפשונו.

**כיצד נמצא את הגרדיאנט של  $E$ :** אנו נחפש את הגרדיאנט של  $\log(P)$  כך

$$\mathcal{L}_{mse} = E_{\mathbf{x} \sim p(\mathbf{x})} [\|\mathcal{F}_\theta(\mathbf{x}) - \nabla_{\mathbf{x}} \log p(\mathbf{x})\|_2^2]$$

נשים לב שזה יתרון, כי גרדיאנט של לוג הוא מנורמל, ולרשתות נוירונים יותר קל לעבוד כך. **הבעיה העיקרית בנוסחה הנ"ל:** אין לנו את  $p$ . **לכן נשתמש בנוסחה הבאה:**

$$\mathcal{L}_{\text{matching}} = E_{\mathbf{x} \sim p(\mathbf{x})} \left[ \text{tr}(\nabla_{\mathbf{x}} \mathcal{F}_\theta(\mathbf{x})) + \frac{1}{2} \|\mathcal{F}_\theta(\mathbf{x})\|_2^2 \right]$$

כך אנו לא צריכים ללמוד את  $p$ .

**מימוש:** נדגום  $x$ -ים מתוך ההתפלגות  $p$ , ונרצה למזער את הפונקציה הזאת על הרשת שלנו. איפה שתהיה התפלגות גבוהה של  $p$  שם אנו נרצה למזער את ה  $divergence$  על מקומות שיותר ריאלי שהם שייכים להתפלגות.

**מה זה  $divergence$ :** מקומות שהרבה חיצים (הצבעות של הגרדיאנט) נכנסים אליהם יוגדרו שליליים. מקומות שהרבה יוצאים מהם יוגדרו חיוביים, ומקומות שיש איזון בין מספר החיצים היוצאים והנכנסים יוגדרו כ 0. אנו נשאף למזער את ה  $divergence$ , כלומר שמספר החיצים שיכנסו יהיה גבוה. כך אנו נעודד התכנסות להתפלגות, כי אם יש וקטורים שמצביעים למקום הזה כנראה הוא שייך להתפלגות, לכן נעודד עוד וקטורים להצביע לשם.

ייעול התהליך עם הורדת מימד: נחשב באופן הבא:

$$E_{\mathbf{v} \sim \mathcal{N}(0,1)} E_{\mathbf{x} \sim p(\mathbf{x})} \left[ \mathbf{v}^T \nabla_{\mathbf{x}} \mathcal{F}_{\theta}(\mathbf{x}) \mathbf{v} + \frac{1}{2} \left\| \mathbf{v}^T \mathcal{F}_{\theta}(\mathbf{x}) \right\|_2^2 \right]$$

## 6.1: Denoising Score Matching

ננסה ללמוד את ההתפלגות בשיטה אחרת.

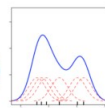
### Denoising Score Matching

[A connection between score matching and denoising autoencoders]

If we corrupt an image by some noise, we get a local distribution  $q_{\sigma}(\tilde{\mathbf{x}} | \mathbf{x})$

It is shown that the following loss can also be used to learn the score

$$E_{q_{\sigma}(\tilde{\mathbf{x}}|\mathbf{x})} E_{\mathbf{x} \sim p(\mathbf{x})} \left[ \left\| \mathcal{F}_{\theta}(\tilde{\mathbf{x}}) - \nabla_{\tilde{\mathbf{x}}} \log q_{\sigma}(\tilde{\mathbf{x}} | \mathbf{x}) \right\|_2^2 \right] \quad k_n = \sum_{i=1}^n \phi\left(\frac{\mathbf{x}-\mathbf{x}_i}{h_n}\right)$$



Parzen Windows

Under the assumption that  $\mathcal{F}_{\theta}(\mathbf{x}) = \nabla_{\mathbf{x}} \log q_{\sigma}(\mathbf{x}) \approx \nabla_{\mathbf{x}} \log p(\mathbf{x})$

In case of a Gaussian noise,  $q_{\sigma}(\tilde{\mathbf{x}} | \mathbf{x}) = \mathcal{N}(\tilde{\mathbf{x}} | \mathbf{x}, \sigma^2 \mathbf{I})$ . we get  $\nabla_{\tilde{\mathbf{x}}} \log q_{\sigma}(\tilde{\mathbf{x}} | \mathbf{x}) = \frac{\tilde{\mathbf{x}} - \mathbf{x}}{\sigma^2}$

Giving the following loss

$$l(\theta; \sigma) = E_{q_{\sigma}(\tilde{\mathbf{x}}|\mathbf{x})} E_{\mathbf{x} \sim p(\mathbf{x})} \left[ \left\| \mathcal{F}_{\theta}(\tilde{\mathbf{x}}) + \frac{\tilde{\mathbf{x}} - \mathbf{x}}{\sigma^2} \right\|_2^2 \right]$$

## 6.2: Diffusion Models

**הרעיון:** מודל של תהליך סטוכסטי, שבכל שלב יש לנו התפלגות נורמלית שמבטאת את ההתפלגות המותנית משלב  $i$  לשלב  $i + 1$ .

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I})$$

$$q(\mathbf{x}_{0:T}) = q(\mathbf{x}_0) \prod_{t=1}^T q(\mathbf{x}_t | \mathbf{x}_{t-1})$$

כאשר השוויון הראשון הוא ההסתברות לעבור מהמצב הנוכחי למצב הבא.

נוכל ללמוד מהסוף להתחלה באופן הבא:

$$p_{\theta}(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t) \quad p(\mathbf{x}_T) = \pi(\mathbf{x}_T)$$

וההסתברות למעבר בין מצבים היא:

$$p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_{\theta}(\mathbf{x}_t, t), \Sigma_{\theta}(\mathbf{x}_t, t))$$

**נוטציות:**  $q$  היא ההתפלגות האמיתית, שיש בעולם. ואנו רויס לקרב אותה באמצעות רשת נוירונים. נסמן ב  $p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)$  את הקירוב של חזרה אחורה.

**אימון הרשת:** כאשר אנו הולכים קדימה התהליך הוא יחסית פשוט, כי אנחנו מתקדמים לכיוון גאוסיאן ומוסיפים רעש בכל שלב, נקודה הופכת לגאוסיאן, וזה דבר שאפשר לצפות. אך בכיוון ההפוך, אנחנו לא יודעים מה המקור ממנו באנו, כי לכל נקודת טשטוש (נקודה שהפכה לגאוסיאן) יכולות להיות כמה מקורות. אנו רוצים שבכל שלב הרשת תפלוט לנו את התוחלת וסטיית התקן.

נקרב את הפתרון על ידי התניה בנוסף על  $X_0$ :

$$q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\mu}(\mathbf{x}_t, \mathbf{x}_0), \tilde{\beta}_t \mathbf{I})$$

**כיצד נחשב את הצעד אחורה:** אנו יודעים את  $q(\mathbf{x}_{t-1} | \mathbf{x}_0)$  ואת  $q(\mathbf{x}_t | \mathbf{x}_0)$ , וגם את  $q(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{x}_0)$ . בעזרת חוק בייס אנו יכולים להסיק את  $q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0)$ .

**פונקציית הלוס שנרצה למזער:**

$$L_t = D_{\text{KL}}(q(\mathbf{x}_t | \mathbf{x}_{t\psi 1}, \mathbf{x}_0) || p_{\theta}(\mathbf{x}_t | \mathbf{x}_{t+1})) \text{ for } 1 \leq t \leq T - 1$$

## 7 Neural Style Transfer

### 7.1 טקסטורה:

**מה הרשת עושה:** נשתמש ברשתות קונבולוציה כדי לייצר תמונות גדולות מדוגמה קטנה - פיתוח טקסטורה.

**הרעיון:** נקח תמונה מורעשת, וננקה אותה עד שהיא תזכיר את התמונה המקורית. נשווה בין התמונות בעזרת רשת VGG שתתן לנו את הפיצרים שנמצאים בתמונה.

**מטריצת גראהם:** עבור כל שני ערוצים בשכבה  $l$  מסויימת, נבדוק את הקורלציה בניהם ונשמור במטריצה  $G$

$$G_{i,j}^l = \frac{1}{C_l H_l W_l} \sum_{h=1}^{H_l} \sum_{w=1}^{W_l} F_{h,w,i} \cdot F_{h,w,j}$$



את הלוס נחשב כך: בעזרת חישוב ההפרשים בין מטריצות גראהם

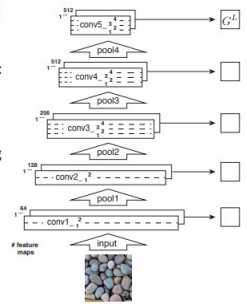
$$L_{style}(x, G(z)) = \sum_{i=1}^5 w_i E_i \quad E_l = \|G^l - A^l\|_F^2 = \sum_{i,j} (G_{i,j}^l - A_{i,j}^l)^2$$

סיכום התהליך:

## Neural Texture Synthesis

1. Pretrain CNN on ImageNet (VGG-19)
2. Feed a texture, record activations
3. Compute Gram matrix for each layer:  

$$G_{ij}^l = \sum_k F_{ik}^l F_{jk}^l$$
4. Initialize image from random noise
5. Feed image through CNN, computing Gram matrices
6. Compute loss as the sum of L2 distances between Gram matrices
7. Back-propagate to get a gradient for the input image
8. Update input image
9. Repeat 5-8 until convergence



## 7.2 היפר ראליסטיקה:

אנו רוצים לקחת תמונה מסויימת ולצייר אותה בסגנון אחר. לדוגמה - להעביר תמונה שצילמנו לציור סטייל ואן-גוך.

הרעיון: נריץ תמונה על רשת VGG, ונשווה את מפות האקטיבציה

### Ingredient 2: Feature Inversion

Given a feature map for an image, find a new image such that:

- Its features are similar to the given features
- It "looks natural" (image prior regularization)

$$\mathbf{x}^* = \underset{\mathbf{x} \in \mathbb{R}^{H \times W \times C}}{\operatorname{argmin}} \ell(\Phi(\mathbf{x}), \Phi_0) + \lambda \mathcal{R}(\mathbf{x})$$

Given feature vector  
Features of new image

$$\ell(\Phi(\mathbf{x}), \Phi_0) = \|\Phi(\mathbf{x}) - \Phi_0\|^2$$

$$\mathcal{R}_{TV^\beta}(\mathbf{x}) = \sum_{i,j} \left( (x_{i,j+1} - x_{ij})^2 + (x_{i+1,j} - x_{ij})^2 \right)^{\frac{\beta}{2}}$$

Total Variation regularizer  
(encourages spatial smoothness)

**מימוש:** נכניס תמונה שלנו ותמונה שנרה את הסטייל שלה, לאחר מכן יהיו לנו אוסף של מטריצות גראהם. נעשה תהליך אופטימיזציה כמו בפעם קודמת.

2.00

# Neural Style Transfer

1. Pretrain CNN
2. Compute features for content image
3. Compute Gram matrices for style image
4. Randomly initialize new image
5. Forward new image through CNN
6. Compute style loss (L2 distance between Gram matrices) and content loss (L2 distance between features)
7. Loss is weighted sum of style and content losses
8. Backprop to image
9. Take a gradient step
10. GOTO 5

