# Practical Part – Bar Rousso I.D. 203765698

## 1. Auto-Encoding

### A) <u>Comparing the SAME encoder/decoder architecture, while changing the latent space dimension</u>
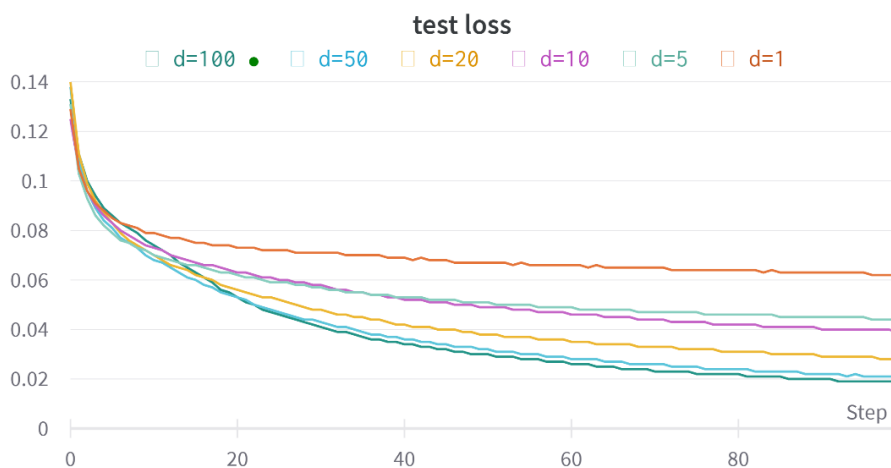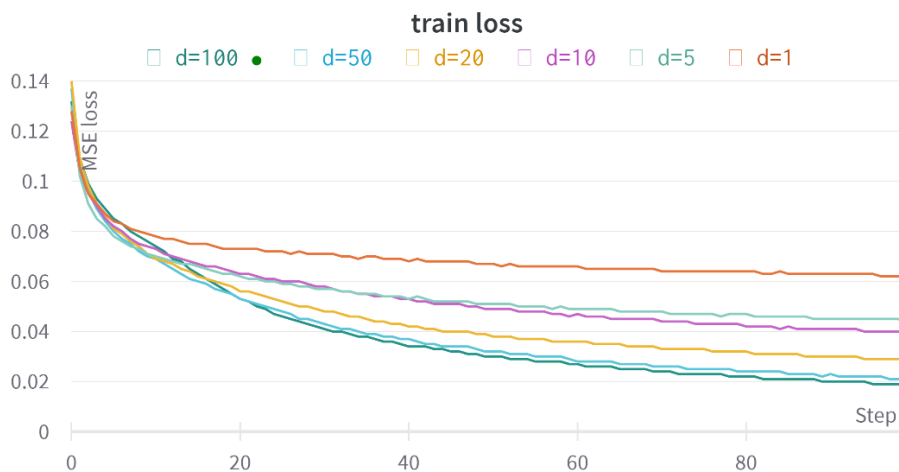
**Encoder network structure: (bais included)**

- Convolution layer: 1 input channel, 18 output channels, 5X5 kernel size, stride = 2
- 2D Batch norm
- RELU
- Convolution layer: 18 input channels, 36 output channels, 3X3 kernel size, stride = 2
- 2D Batch norm
- RELU
- Fully connected layer: 900 input feathers, $d$ output features
- RELU

**Encoder network structure: (bais included)**

- Fully connected layer: $d$ input features, 900 output features
- RELU
- Transpose convolution layer: 36 input channels, 18 output channels, 3X3 kernel size, stride = 2, out padding = 1
- 2D Batch norm
- RELU
- Transpose convolution layer: 18 input channels, 1 output channel, 5X5 kernel size, stride = 2, out padding = 1
- Sigmoid

**Results:**

| Latent dimension | Number of parameters | Train loss | Test loss |
|---|---|---|---|
| 1 | 15482 | 0.062 | 0.062 |
| 5 | 22686 | 0.044 | 0.044 |
| 10 | 31691 | 0.04 | 0.038 |
| 20 | 49701 | 0.029 | 0.028 |
| 50 | 103731 | 0.021 | 0.021 |
| 100 | 193781 | 0.019 | 0.019 |

**Conclusions:**

We can see that as we increase the latent space dimension, we get **lower** train and test losses. This can be explained as bigger latent space dimension enable each encoded vector to captures more details of its original image.

## B) Comparing encoder/decoder architectures with different number of layers, while FIXING the latent space dimension

One layer architecture:

**Encoder network structure: (bais included)**

- Convolution layer: 1 input channel, 18 output channels, 5X5 kernel size, stride = 4
- 2D Batch norm
- RELU
- Fully connected layer: 648 input feathers, $d$ output features
- RELU

**Encoder network structure: (bais included)**

- Fully connected layer: $d$ input features, 648 output features
- Transpose convolution layer: 18 input channels, 1 output channel, 5X5 kernel size, stride = 4, out_padding = 3
- Sigmoid

Total number of parameters: 27543

Two layers architecture: The same architecture described in the first section

Three layers architecture:

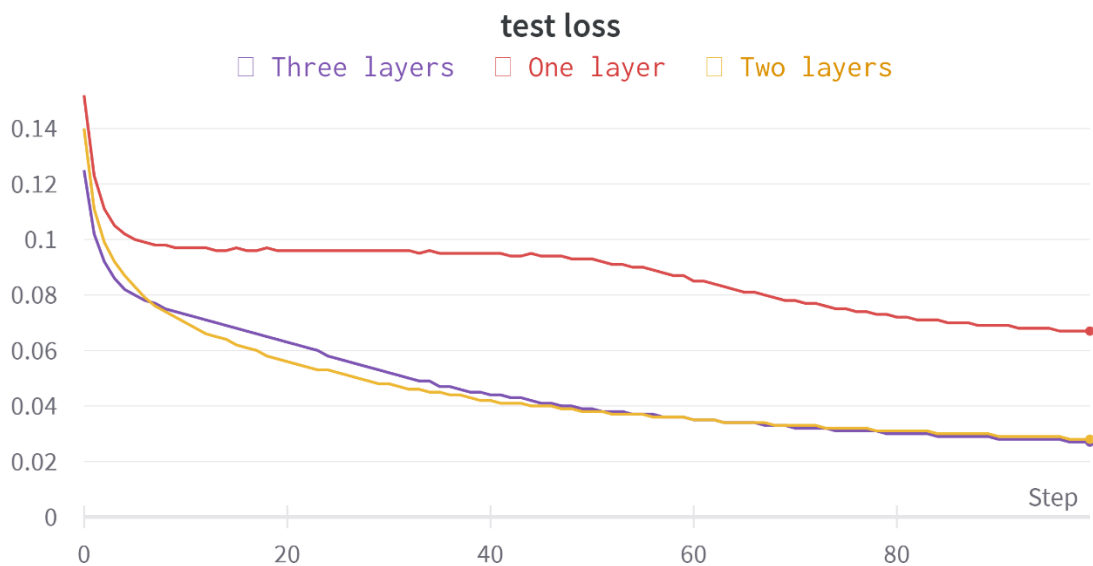**Encoder network structure: (bais included)**

- Convolution layer: 1 input channel, 9 output channels, 3X3 kernel size, stride = 2, padding = 1
- 2D Batch norm
- RELU
- Convolution layer: 9 input channels, 18 output channels, 3X3 kernel size, stride = 2, padding = 1
- 2D Batch norm
- RELU
- Convolution layer: 18 input channels, 36 output channels, 3X3 kernel size, stride = 2, padding = 1
- 2D Batch norm
- RELU
- Fully connected layer: 324 input feathers, $d$ output features
- RELU

**Encoder network structure: (bais included)**

- Fully connected layer: $d$ input features, 324 output features
- RELU
- Transpose convolution layer: 36 input channels, 18 output channels, 3X3 kernel size, stride = 2, out_padding = 1
- 2D Batch norm
- RELU
- Transpose convolution layer: 18 input channels, 9 output channels, 3X3 kernel size, stride = 2, padding = 1, out_padding = 1
- 2D Batch norm
- RELU
- Transpose convolution layer: 9 input channels, 1 output channel, 3X3 kernel size, stride = 2, padding = 1, out_padding = 1
- Sigmoid

Total number of parameters: 28317

**Results:**



**train loss**

Three layers    One layer    Two layers



**test loss**

Three layers    One layer    Two layers

**Conclusions:**

We can see that as we add more convolutions layers, we got a better module in terms of:

(1) Lower final train and test losses results

(2) Faster convergence to a module with loss train and test losses

(3) Using less parameters for better performance
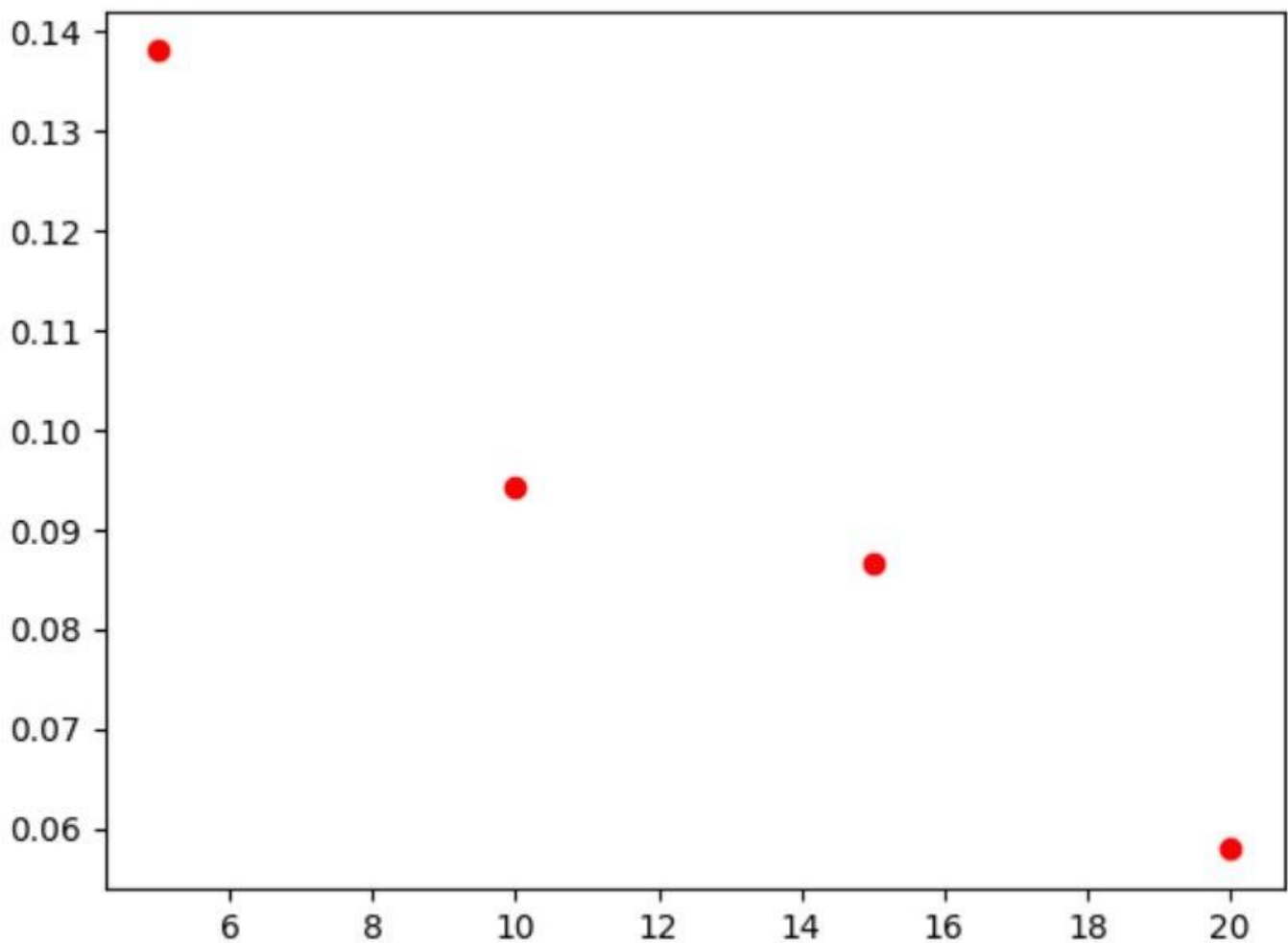
## 2. **Interpolation**

Comparing results generated from different latent space dimensions 20 VS 100:

| Interpolation between | latent space dimension | Results |
|---|---|---|
| 3 -> 5 | 20 |  |
| | 100 |  |
| 5 -> 8 | 20 |  |
| | 100 |  |
| 0 -> 4 | 20 |  |
| | 100 |  |
| 2 -> 9 | 20 |  |
| | 100 |  |

**Conclusions:** We can see that the bigger latent space dimension we use, the decoder produce more "sharper" images that look more realistic. This is because a bigger latent space dimension enables us to encode more details of the image.

### 3. **Decorrelation**

Plot of the MEAN Pearson correlation between all couples of coordinates in the latent space (in absolute values) as a function of the latent space dimension:



**Conclusions:**

As we can see in the plot, as the we increase the latent space dimension, we get a smaller correlation between different coordinates. This is because a larger latent space dimension enables us to better represent the images by capturing **more variations** of the image, resulting in more coordinates with **wicker correlation**.

**Note:** Some pairs of the Pearson correlation resulted in NAN. To overcome this, I replaced them with zeros to not influence the mean.

## 4. **Transfer Learning**

**Trained Encoder network structure: (bais included, latent space dimension = 20)**

- Convolution layer: 1 input channel, 9 output channels, 3X3 kernel size, stride = 2, padding = 1
- 2D Batch norm
- RELU
- Convolution layer: 9 input channels, 18 output channels, 3X3 kernel size, stride = 2, padding = 1
- 2D Batch norm
- RELU
- Convolution layer: 18 input channels, 36 output channels, 3X3 kernel size, stride = 2, padding = 1
- 2D Batch norm
- RELU
- Fully connected layer: 324 input feathers, 20 output features
- RELU

**MLP network structure: (bais included)**

- Fully connected layer: 20 input feathers, 50 output features
- RELU
- Fully connected layer: 50 input feathers, 100 output features
- RELU
- Fully connected layer: 100 input feathers, 10 output features
- SoftMax

The **Classification** module concrete both encoder and MLP, and was trained with respect to Cross Entropy loss function.

The classification model was trained according to two scenarios:

1. Only MLP wights were updated
2. Both MLP and Encoder wights were updated

**Results:**

Scenario 1: `train loss=2.303, test loss=2.304`

Scenario 2: `train loss=2.300, test loss=2.300`

**Comment:**

Unfortunately, I got losses values that is not make sense. I didn't succeed to debug the problem.

I know that the second scenario supposed to have better results, since we allow to more parameters to be updated.

Also, as I mentioned above, the bigger the latent space dimension, the better the latent vectors as they capturing more details of the input images.

Thus, the MLP module will do a better job in classify the vectors correctly as they contain more data.

א) כ"ג: הרכבה של פונקציות ליניאריות הינה פונקציה ליניארית שנוכיח.

הוכחה: יהיו $f: B \to C$   $g: A \to B$ שני פונקציות ליניאריות מעל מרחבים וקטוריים $A, B, C$ מעל אותו שדה $R$.

נראה כי $f \circ g : A \to C$ והן גם ליניאריות נראה.

(I) אדיטיביות: יהיו $a_1, a_2 \in A$ אז

$$f \circ g(a_1 + a_2) \underset{\text{אדיטיביות של } g}{=} f(g(a_1) + g(a_2)) \underset{\text{אדיטיביות של } f}{=} f \circ g(a_1) + f \circ g(a_2) \,\|$$

(II) הומוגניות: יהי $a \in A$  ו-  $\alpha \in R$

$$f \circ g(\alpha \cdot a) \underset{\text{הומוגניות של } g}{=} f(\alpha \cdot g(a)) \underset{\text{הומוגניות של } f}{=} \alpha \, f \circ g(a)$$

אז הראנו נכון כ $f \circ g$ פונקציה ליניארית שנרצה.

ב) כ"ג: הרכבה של פונקציות אפיניות הינה פונקציה אפינית שנוכיח.

הוכחה: יהיו $g: A \to B$, $f: B \to C$ שני פונקציות אפיניות מעל מרחבים וקטוריים $A, B, C$ מעל אותו שדה $R$.

על הגדרה הפונקציות האפיניות קיימים $b^* \in B$, $c^* \in C$ והמטריצות האפיניות $M_g : A \to B$ , $M_f : B \to C$ כך ש:

$g(a) = M_g(a) + b^*$ לכל $a \in A$

$f(b) = M_f(b) + c^*$ לכל $b \in B$

ולכן לכל $a \in A$:

$$f \circ g(a) = f(M_g(a) + b^*) = M_f(M_g(a) + b^*) + c^*$$
$$\underset{\text{אדיטיביות של } M_f}{=} M_f \circ M_g(a) + M_f(b^*) + c^* \,\|$$

אנו $M_f \circ M_g$ פונקציה ליניארית (מכפלת מטריצות) וכן $M_f(b^*) + c^*$ הוא וקטור קבוע לכן נכון כי $f \circ g$ פונקציה אפינית שנרצה.

2) א) סביר ... האימון ... בכלל ... $\theta^{n+1} = \theta^n - \alpha \nabla f_{\theta^n}(x)$

אז האם ... תמיד?

... $f_{\theta^n}$ ... $f_{\theta^n}$. ... $f_{\theta^n}$ ...

- אפשר ... על ...

$\left| f_{\theta^n}(\theta^n) - f_{\theta^n}(\theta^{n-1}) \right| < \varepsilon$ - ...

...

ב) ... $x_0$ ...

$f$ ...:

$$f(x) = f(x_0) + \nabla f(x_0) \cdot (x - x_0) + (x - x_0)^T \cdot H(x_0) \cdot (x - x_0) + o(\|x - x_0\|^3)$$

כאשר $H(x_0)$ ... של $f$ בנקודה $x_0$

... $x_0$ ...

... $f$.

(I) ... אם $H(x_0)$ ... $x_0$ ... $f$.

אבל, אם $H(x_0)$ ... $x_0 \neq x$ ...

$(x - x_0)^T \cdot H(x_0) \cdot (x - x_0) > 0$

rtl
אזן נובע (וכן) 6 דרוש כי נוכל לקזל לכל שבל

א~דר"ק $X_0 \neq X \in (X_0 - \epsilon, X_0 + \epsilon)$

$$(X - X_0)^T H(X_0)(X - X_0) + o(||X - X_0||^3) \geq 0$$

$$f(X) = f(X_0) + \underbrace{\nabla f(X_0)(X - X_0)}_{0} + \underbrace{(X - X_0)^T \cdot H(X_0) \cdot (X - X_0) + o(||X - X_0||^3)}_{\underset{0}{\lor}} \overset{\text{של}}{\geq}$$

$$\geq f(X_0)$$

(וכן) כי $X_0$ זוהי כך נק' מינימום מקומי של $f$.

II) כאן $H(X_0)$ מטריצה מוגדרת שלילית נקבל שלכל

$X_0$ כך נק' מקסימום מקומי של $f$.

לכן, נניח ש־ $H(X_0)$ מטריצה מוגדרת שלילית, אזי לכל $X \neq X_0$

$$(X - X_0)^T H(X_0) \cdot (X - X_0) < 0$$

ובאופן דומה, (וכן) כי דרוש 6 נוכל לקזל לכל שבל

א~דר"ק $X_0 \neq X \in (X_0 - \epsilon, X_0 + \epsilon)$

$$(X - X_0)^T H(X_0)(X - X_0) + o(||X - X_0||^3) \leq 0$$

$$f(X) = f(X_0) + \underbrace{\nabla f(X_0)(X - X_0)}_{0} + \underbrace{(X - X_0)^T \cdot H(X_0) \cdot (X - X_0) + o(||X - X_0||^3)}_{\underset{0}{\land}} \overset{\text{של}}{\leq}$$

$$\leq f(X_0)$$

(וכן) כי $X_0$ זוהי נק' מקסימום מקומי של $f$.

נניח כי... מכיוון... שערך... תחום... אליו (0°-360°)

שערך... מן סופר... כך... מחבר... על... מ... של...

כך אותה... Loss(2,360)=Loss(0,2)

## Pseudo-code:

```
def loss ( pred_deg, true_deg):


    rad_pred_deg = 2π·(pred_deg/180)
    rad_true_deg = 2π·(true_deg/180)


    while rad_pred_deg < 0:
        rad_pred_deg += 2π

    while rad_true_deg < 0:
        rad_true_deg += 2π


    difference = | rad_true_deg - rad_pred_deg |

    result = min { difference, 2π - difference }
    return result
```

a) $\dfrac{\partial}{\partial x} f(x+y, 2x, z)$

נגדיר את הפונקציה $g: \mathbb{R}^3 \to \mathbb{R}^3$ ש' $g(x,y,z) = (x+y, 2x, z)$

$\dfrac{\partial}{\partial x} f \circ g(x,y,z)$   אל נמצא את הגזירות ואז נשים לבל כל נעשה ר

$D_{f \circ g}(x,y,z) = D_f(g(x,y,z)) \cdot D_g(x,y,z) =$

$= D_f(x+y, 2x, z) \cdot \begin{pmatrix} 1 & 1 & 0 \\ 2 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix} = \nabla_f(x+y, 2x, z)^T \begin{pmatrix} 1 & 1 & 0 \\ 2 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix}$

עכשיו נוכל לפי ל:

$\dfrac{\partial}{\partial x} f(x+y, 2x, z) = \left( D_{f \circ g}(x,y,z) \right)_1 = \nabla_f(x+y, 2x, z) \cdot \begin{pmatrix} 1 \\ 2 \\ 0 \end{pmatrix}$

↑
כרכיב הראשון

b) $f_1(f_2(\ldots f_n(x)))$

עכ' כל בכלל:

$D_{f_1 \circ \ldots \circ f_n}(x) = D_{f_1 \circ \ldots \circ f_{n-1}}(f_n(x)) \cdot D_{f_n}(x) =$

$= D_{f_1 \circ \ldots \circ f_{n-2}}(f_{n-1} \circ f_n(x)) \cdot D_{f_{n-1}}(f_n(x)) \cdot D_{f_n}(x) =$

$= D_{f_1 \circ \ldots \circ f_{n-3}}(f_{n-2} \circ f_{n-1} \circ f_n(x)) \cdot D_{f_{n-2}}(f_{n-1} \circ f_n(x)) \cdot D_{f_{n-1}}(f_n(x)) \cdot D_{f_n}(x) =$

$= \ldots = D_{f_1}(f_2 \circ \ldots \circ f_n(x)) \cdot D_{f_2}(f_3 \circ \ldots \circ f_n(x)) \cdot \ldots \cdot D_{f_n}(x) =$

$= \displaystyle\prod_{i=1}^{n} D_{f_i}(f_{i+1} \circ \ldots \circ f_n(x))$

c) $f_1(x, f_2(x, f_3(\cdots f_{n-1}(x, f_n(x)))$

$f_n : \mathbb{R} \to \mathbb{R}$ ‏כאשר‏ ‏נניח‏

$f_i : \mathbb{R}^2 \to \mathbb{R}$ ‏ $1 \leq i \leq n-1$‏ ‏ולכל‏

$f_{n-1}(x, f_n(x)) = f_{n-1} \circ g_n (x)$     ‏נגדיר‏ $g_n(x) = (x, f_n(x))$

‏אזי:‏

$\left( f_{n-1} \circ g_n (x) \right)' = D_{f_{n-1}}(g_n(x)) \cdot D_{g_n}(x) =$

$= \nabla f_{n-1}(x, f_n(x))^T \cdot \begin{pmatrix} 1 \\ f_n'(x) \end{pmatrix}$

‏נגדיר‏ $g_{n-1}(x) = (x, f_{n-1} \circ g_n(x))$ :

$f_{n-2}(x, f_{n-1}(x, f_n(x))) = f_{n-2} \circ g_{n-1}(x)$

‏אזי:‏

$\left( f_{n-2} \circ g_{n-1}(x) \right)' = D_{f_{n-2}}(g_{n-1}(x)) \cdot D_{g_{n-1}}(x) =$

$= \nabla f_{n-2}(x, f_{n-1} \circ g_n(x))^T \cdot \begin{pmatrix} 1 \\ (f_{n-1} \circ g_n(x))' \end{pmatrix}$

‏נגדיר‏ $g_{n-2}(x) = (x, f_{n-2} \circ g_{n-1}(x))$

$f_{n-3}(x, f_{n-2}(x, f_{n-1}(x, f_n(x)))) = f_{n-3} \circ g_{n-2}(x)$

$\left( f_{n-3} \circ g_{n-2}(x) \right)' = D_{f_{n-3}}(g_{n-2}(x)) \cdot D_{g_{n-2}}(x) =$    ‏אזי‏

$= \nabla f_{n-3}(x, f_{n-2} \circ g_{n-1}(x))^T \cdot \begin{pmatrix} 1 \\ (f_{n-2} \circ g_{n-1}(x))' \end{pmatrix}$

$$g_2(x) = (X, f_2 \circ g_3(x))$$ נקבל כי

$$f_1(X, f_2(X, f_3(\cdots f_{n-1}(X, f_n(X)))) = f_1 \circ g_2(x)$$

לכן

$$(f_1 \circ g_2(x))' = D_{f_1}(g_2(x)) \cdot D_{g_2}(x) =$$

$$= D_{f_1}(X, f_2 \circ g_3(x))^T \cdot \begin{pmatrix} 1 \\ (f_2 \circ g_3(x))' \end{pmatrix} \quad //$$

נוכל לכ לכל $3 \leq i \leq n$:

$$f_{i-1} \circ g_i(x) = \begin{cases} f_{i-1}(X, f_i \circ g_{i+1}(x)) & \text{if } i \leq n-1 \\ f_{n-1}(X, f_n(x)) & \text{if } i = n \end{cases}$$

$$(f_{i-1} \circ g_i(x))' = \begin{cases} D_{f_1}(X, f_i \circ g_{i+1}(x))^T \cdot \begin{pmatrix} 1 \\ (f_i \circ g_{i+1}(x))' \end{pmatrix} & \text{if } i \leq n-1 \\ D_{f_{n-1}}(X, f_n(x))^T \cdot \begin{pmatrix} 1 \\ f_n'(x) \end{pmatrix} & \text{if } i = n \end{cases}$$

d) $f(x + g(x + h(x)))$

$k_1(x) = x + h(x)$

$k_2(x) = x + g \circ k_1(x)$

$$D_{f \circ k_2}(x) = D_f(k_2(x)) \cdot D_{k_2}(x) =$$

$$= D_f(x + g \circ k_1(x)) \cdot (I_n + D_{g \circ k_1}(x)) =$$

$$= D_f(x + g(x + h(x))) \cdot (I_n + D_g(k_1(x)) \cdot D_{k_1}(x)) =$$

$$= D_f(x + g(x + h(x))) \cdot (I_n + D_g(x + h(x)) \cdot (I_n + D_h(x))) =$$

$$= D_f(x + g(x + h(x))) \cdot (I_n + D_g(x + h(x)) + D_g(x + h(x)) \cdot D_h(x)) //$$

$$\left(f \circ k_2(x)\right)' = f'(x + g(x + h(x))) \cdot (1 + g'(x + h(x)) + g'(x + h(x)) \cdot h'(x)) //$$