

# Practical part – Bar Rouso 203765698

## Q1:

Base model (what I started with):

**network structure: (bais included)**

- Convolution layer: 3 input channels, 6 output channels, 5x5 kernel size
- RELU
- Max pooling: 2x2 Window size
- Convolution layer: 6 input channels, 16 output channels, 5x5 kernel size
- RELU
- Max pooling: 2x2 Window size
- Fully connected layer: 400 input features, 10 output features

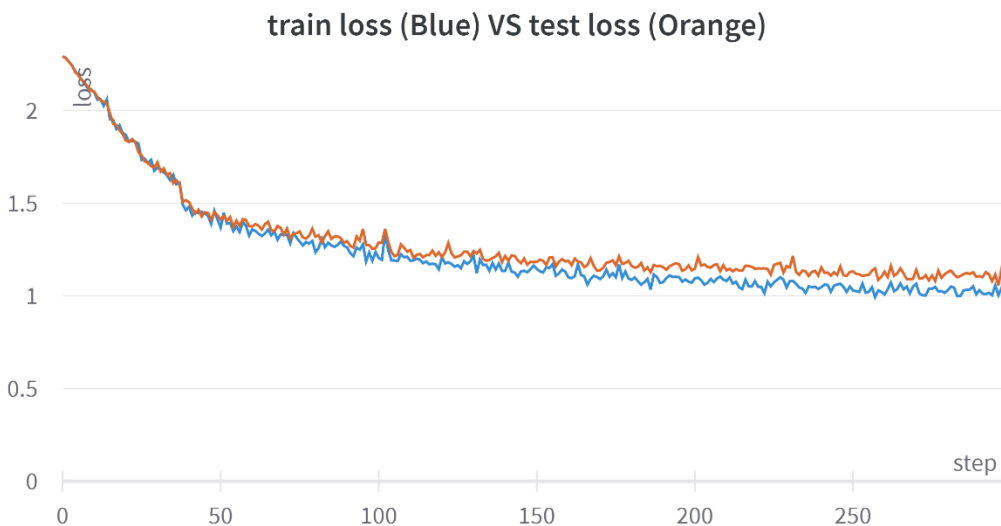
**Total learnable parameters:**

$$(5 \cdot 5 \cdot 3 + 1) \cdot 6 + (5 \cdot 5 \cdot 6 + 1) \cdot 16 + (16 \cdot 5 \cdot 5 + 1) \cdot 10 = 6882$$

**Learning rate:** 0.001

**batch size:** 32

**Number of epochs:** 10



**Comment:** In this part I only first use the original network from the tutorial and observed the results.

### Overfitting model:

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

- Convolution layer: 3 input channels, 75 output channels, 5x5 kernel size
- RELU
- Max pooling: 2x2 Window size
- Convolution layer: 75 input channels, 150 output channels, 5x5 kernel size
- RELU
- Max pooling: 2x2 Window size
- Fully connected layer: 3750 input features, 10 output features

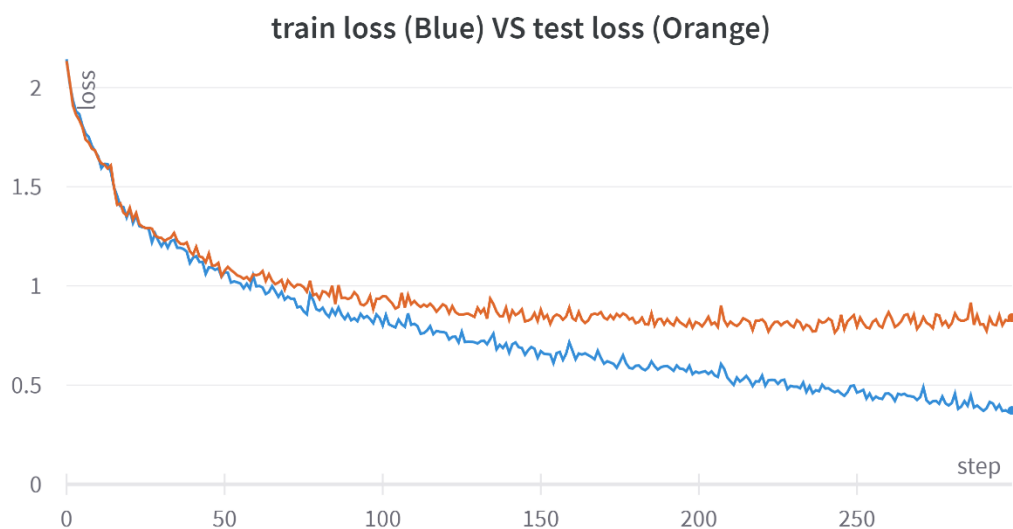
#### **Total learnable parameters:**

$$(5 \cdot 5 \cdot 3 + 1) \cdot 75 + (5 \cdot 5 \cdot 75 + 1) \cdot 150 + (150 \cdot 5 \cdot 5 + 1) \cdot 10 = 324,610$$

**Learning rate:** 0.001

**batch size:** 32

**Number of epochs:** 20



**Comment:** In order to achieve overfitted model I did:

- Increased the number of learnable parameters to create more complex hypothesis class.
- Increase the number of epochs to make the model "memorize" the training set, by repeatedly going through the training set and keep adjust the model to better match the training results.

### Underfitting model:

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

- Convolution layer: 3 input channels, 2 output channels, 3x3 kernel size
- RELU
- Max pooling: 2x2 Window size
- Convolution layer: 2 input channels, 2 output channels, 3x3 kernel size
- RELU
- Max pooling: 2x2 Window size
- Fully connected layer: 72 input features, 10 output features

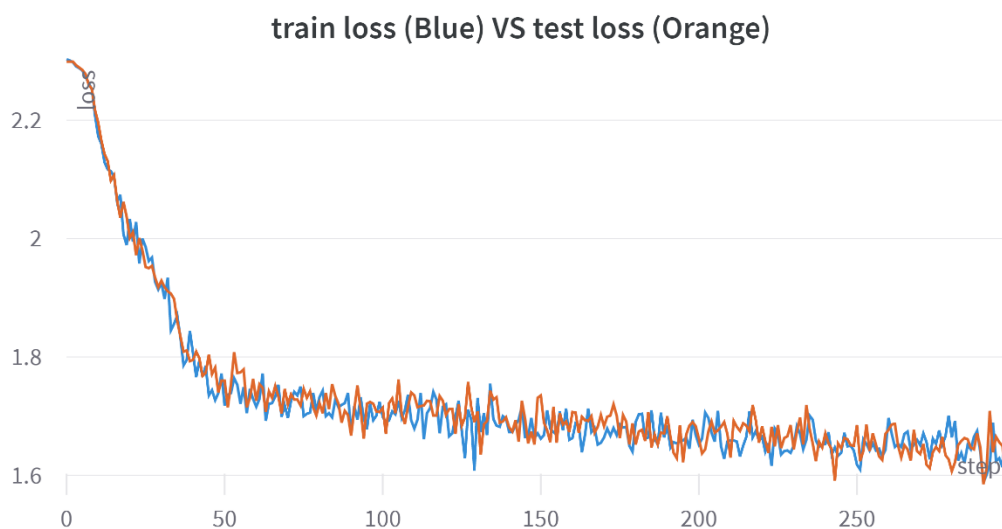
#### **Total learnable parameters:**

$$(3 \cdot 3 \cdot 3 + 1) \cdot 2 + (3 \cdot 3 \cdot 2 + 1) \cdot 2 + (2 \cdot 6 \cdot 6 + 1) \cdot 10 = 824$$

**Learning rate:** 0.01

**batch size:** 32

**Number of epochs:** 5



**Comment:** In order to achieve underfitted model I did:

- Decreased the number of learnable parameters to create more simple hypothesis class.
- Decreased the number of epochs, such that the model didn't go through enough steps to reach to local minimum in SGD process.
- Decreased the learning rate, so the model can miss a local minimum.

Best model (I found):

**network structure: (bais included)**

- Convolution layer: 3 input channels, 18 output channels, 5x5 kernel size
- RELU
- Max pooling: 2x2 Window size
- Convolution layer: 18 input channels, 36 output channels, 5x5 kernel size
- RELU
- Max pooling: 2x2 Window size
- Fully connected layer: 900 input features, 10 output features

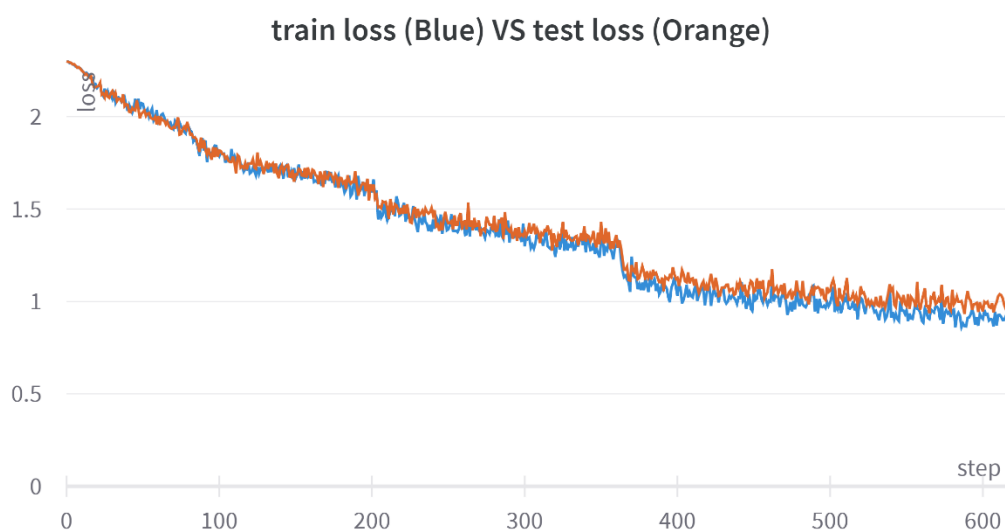
**Total learnable parameters:**

$$(5 \cdot 5 \cdot 3 + 1) \cdot 18 + (5 \cdot 5 \cdot 18 + 1) \cdot 36 + (36 \cdot 5 \cdot 5 + 1) \cdot 10 = 26614$$

**Learning rate:** 0.0005

**batch size:** 16

**Number of epochs:** 10



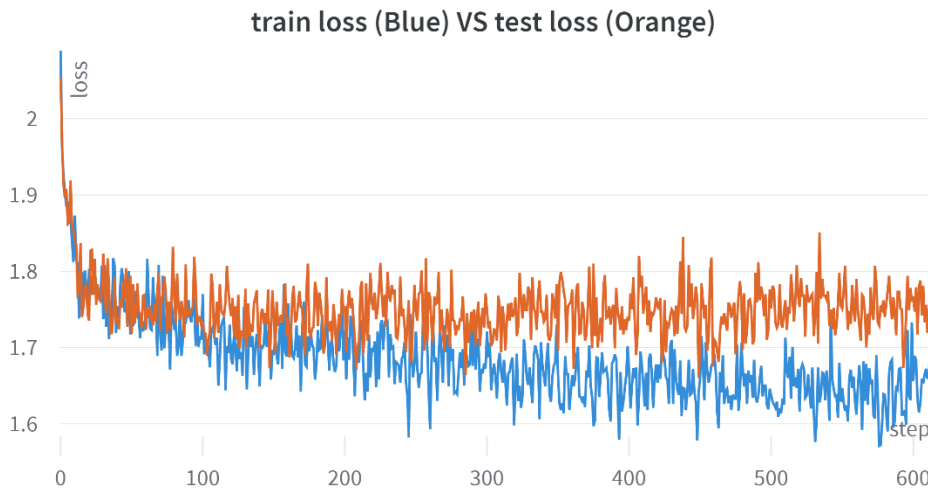
**Comment:** In order to adjust the model to achieve better results I did:

- Used 10 epochs as a trade off regards to the underfitted and overfitted models.
- Used not too many and not too less parameters, regard to the underfitted and overfitted models.
- Decrease learning rate to not miss a local minimum during SGD process.

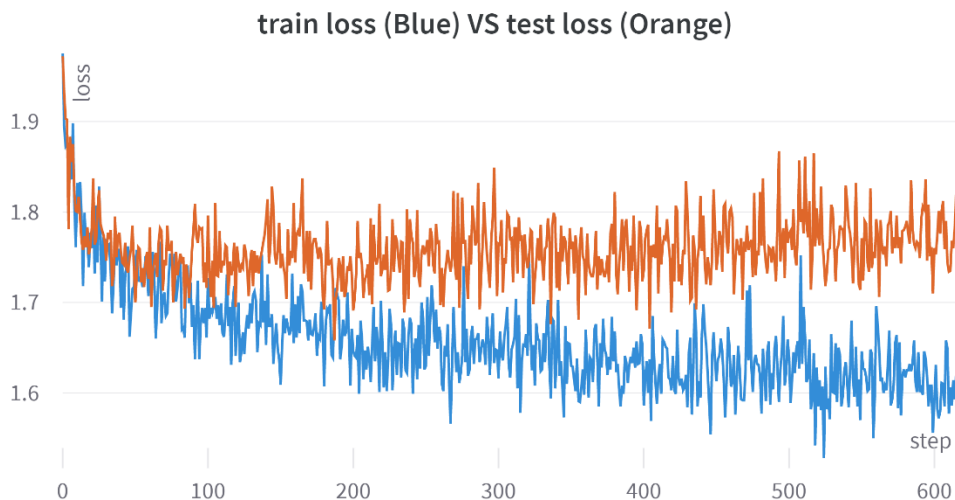
## Q2:

The network's non-linear components are RELU and max pool operators.  
Removing them will result in the following network structure:

- Convolution layer: 3 input channels, 18 output channels, 5x5 kernel size
- Convolution layer: 18 input channels, 36 output channels, 5x5 kernel size
- Fully connected layer: 20736 input features, 10 output features



Increasing the network's size didn't not changed much:



The reason that non-linear components are important:

Supposing that we have a neural network with  $T$  linear layers.

Since each layer  $L_t$  is linear, we can express its operation as a matrix multiplication  $W_t$

Thus, the whole network can be represented by a **single** linear layer, which represented by the matrix:

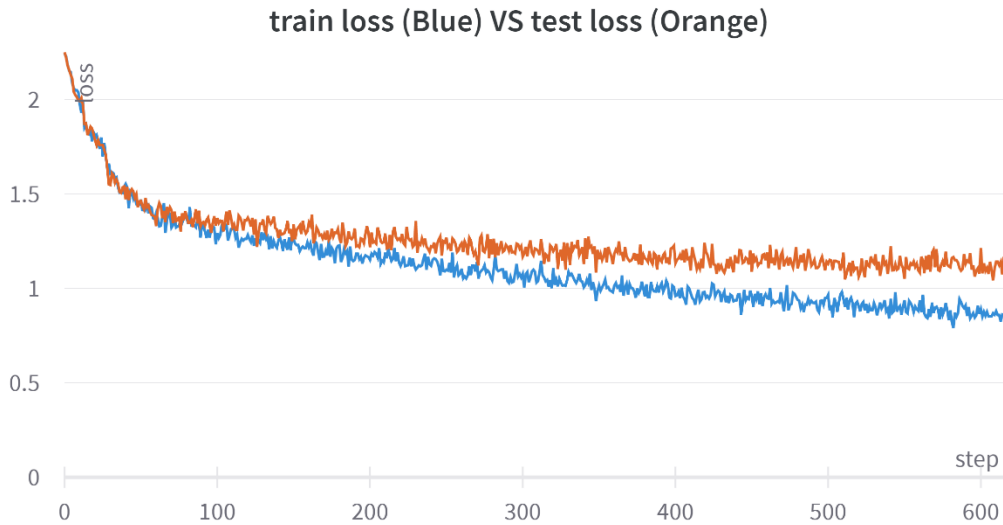
$$W_T \cdot W_{T-1} \cdot \dots \cdot W_1.$$

For that reason, adding only linear layers will not increase the approximation power of a linear neural network at all.  
This is the reason we get underfitting model with poor results.

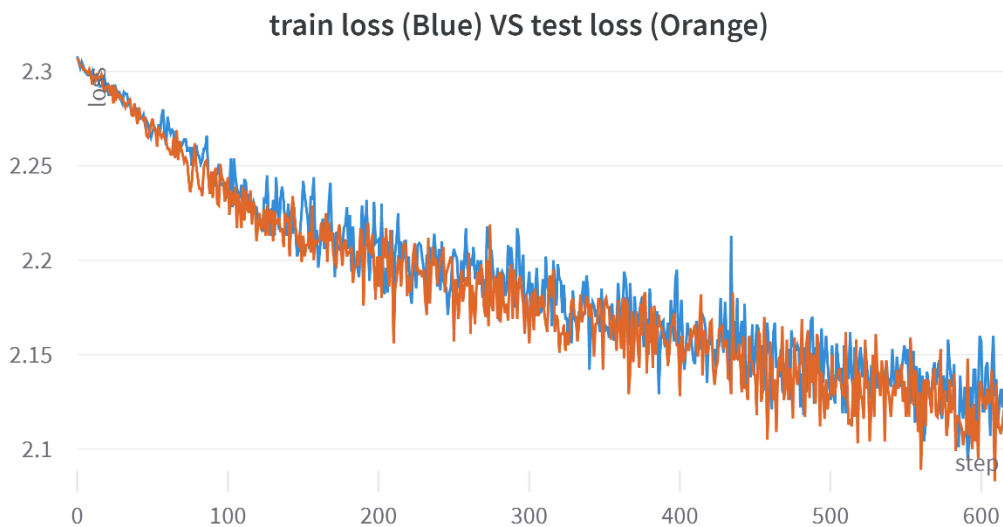
### Q3:

With the current network architecture (from Q1) this are the plots I got for each option:

Option1: A FC layer on the activations of the first conv layer



Option2: A “global average pooling” operator that compute the average of the activations in each channel and then apply FC layer.



We can see that we get better results on the first option, so we choose **option number 1**.

This makes sense, because when we use the global average pooling, we basically only consider only the average value on each channel to make a further decision about the image subject.

But since many images with subjects of different classes can share similar average values on each channel, we can expect to many misclassifications.

However, after I adapted the network architecture with option1, the new network didn't perform better compare to the original network from Q1.

This is makes sense, because in the new network, the receptive field of each neuron in the output layer is **smaller** compare to the original network.

Thus, the new module will be difficult in recognize basic shapes that spread in a large area of the image, which can give a useful information regards to the image's subject classification.



הפך מ'אברהם' - קר לוח 203465698

$$(1) \quad \text{נניח } (Lx_{i+k})_j = (Lx_i)_{j+k} \quad \text{כאשר } x \in \mathbb{R}^n \text{ ו- } L \text{ מטריצה } n \times n$$

אם  $x \in \mathbb{R}^n$  ו-  $L$  מטריצה  $n \times n$  אז  $(Lx)_{i+k} = (Lx)_i + k$  (מודול  $n$ )  
 כלומר,  $(Lx)_{i+k} = (Lx)_i + k$  (מודול  $n$ )  
 כלומר,  $(Lx)_{i+k} = (Lx)_i + k$  (מודול  $n$ )

נניח  $Lx_{i+k} = (Lx_i)_{j+k}$  (מודול  $n$ )

$$(Lx_{i+k})_j = \sum_{a=1}^m (x_{i+k})_{j-a \pmod n} \cdot (L)_{a,j} = \sum_{a=1}^m (x_i)_{j+k-a \pmod n} \cdot (L)_{a,j} = (**)$$

כלומר,  $(Lx_{i+k})_j = (Lx_i)_{j+k}$  (מודול  $n$ )  
 כלומר,  $(Lx_{i+k})_j = (Lx_i)_{j+k}$  (מודול  $n$ )  
 כלומר,  $(Lx_{i+k})_j = (Lx_i)_{j+k}$  (מודול  $n$ )

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

כלומר,  $(Lx_{i+k})_j = (Lx_i)_{j+k}$  (מודול  $n$ )



(2)  $V_{in} \in \mathbb{R}^n$   $FC: \mathbb{R}^n \rightarrow \mathbb{R}^m$   $\rightarrow$   $V_{in}$   $\in \mathbb{R}^n$   $\rightarrow$   $V_{in}$   $\in \mathbb{R}^n$   $\rightarrow$   $V_{in}$   $\in \mathbb{R}^n$

משפט:  $FC$   $\rightarrow$   $W \in \mathbb{R}^{m \times n}$

הוכחה:  $FC$   $\rightarrow$   $W \in \mathbb{R}^{m \times n}$   $\rightarrow$   $W \in \mathbb{R}^{m \times n}$

$W \in \mathbb{R}^{m \times n}$   $\rightarrow$   $W \in \mathbb{R}^{m \times n}$

$FC$   $\rightarrow$   $W \in \mathbb{R}^{m \times n}$   $\rightarrow$   $W \in \mathbb{R}^{m \times n}$

$(W_1, \dots, W_n) = W \in \mathbb{R}^{m \times n}$

$FC$   $\rightarrow$   $W \in \mathbb{R}^{m \times n}$   $\rightarrow$   $W \in \mathbb{R}^{m \times n}$

$FC$   $\rightarrow$   $W \in \mathbb{R}^{m \times n}$   $\rightarrow$   $W \in \mathbb{R}^{m \times n}$

$(FC(V_{in}))_j = (V_{in})_{G^{-1}(j)}$   $\rightarrow$   $(FC(V_{in}))_j = (V_{in})_{G^{-1}(j)}$

$FC$   $\rightarrow$   $W \in \mathbb{R}^{m \times n}$   $\rightarrow$   $W \in \mathbb{R}^{m \times n}$

$FC$   $\rightarrow$   $W \in \mathbb{R}^{m \times n}$   $\rightarrow$   $W \in \mathbb{R}^{m \times n}$

$FC$   $\rightarrow$   $W \in \mathbb{R}^{m \times n}$   $\rightarrow$   $W \in \mathbb{R}^{m \times n}$

$(W_1, \dots, W_n) = W \in \mathbb{R}^{m \times n}$   $\rightarrow$   $(W_1, \dots, W_n) = W \in \mathbb{R}^{m \times n}$

$FC$   $\rightarrow$   $W \in \mathbb{R}^{m \times n}$   $\rightarrow$   $W \in \mathbb{R}^{m \times n}$

$$W_j = W_{G^{-1}(j)}$$

$FC$   $\rightarrow$   $W \in \mathbb{R}^{m \times n}$   $\rightarrow$   $W \in \mathbb{R}^{m \times n}$

$$(W \cdot FC(V_{in}))_j = \langle -W_j, G(V_{in}) \rangle = \sum_{k=1}^n (W_j)_k \cdot (G(V_{in}))_k =$$

$$= \sum_{k=1}^n W_{j \cdot G^{-1}(k)} \cdot (V_{in})_{G^{-1}(k)} = \sum_{k=1}^n W_{j \cdot k} \cdot (V_{in})_k =$$

$$= \langle -W_j, V_{in} \rangle = (W \cdot V_{in})_j$$

$FC$   $\rightarrow$   $W \in \mathbb{R}^{m \times n}$   $\rightarrow$   $W \in \mathbb{R}^{m \times n}$



(3)  $\text{ReLU}$  פונקציה ליניארית (LTI)

(b)  $\text{ReLU}$

links  $\rightarrow$   $V \in \mathbb{R}^n$ ,  $j \leq n$

$$(\text{ReLU}(V))_j = \max\{0, V_j\}$$

מיון  $\text{ReLU}$  אינו פונקציה ליניארית

$$\text{ReLU}\left(\begin{pmatrix} 2 \\ 2 \end{pmatrix}\right) + \text{ReLU}\left(\begin{pmatrix} -1 \\ -1 \end{pmatrix}\right) = \begin{pmatrix} 2 \\ 2 \end{pmatrix} \neq \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \text{ReLU}\left(\begin{pmatrix} 2 \\ 2 \end{pmatrix} + \begin{pmatrix} -1 \\ -1 \end{pmatrix}\right)$$

לכן  $\text{ReLU}$  אינו LTI (לפי הניסוי)

(c) Strided Pooling Layer

שכבה  $M$  בגודל  $m \times n$

$$L(M) = M_{1,n} \text{ (for right)}$$

שכבה  $L$  בגודל  $l \times n$ , LTI

$$M_1 = \begin{pmatrix} 3 & 3 & 3 & 3 \\ 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 \\ 2 & 2 & 2 & 2 \end{pmatrix} \quad M = \begin{pmatrix} 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 \\ 2 & 2 & 2 & 2 \\ 3 & 3 & 3 & 3 \end{pmatrix}$$

$$L(M_1) = \begin{pmatrix} 3 & 3 \\ 1 & 1 \end{pmatrix} \quad L(M) = \begin{pmatrix} 0 & 0 \\ 2 & 2 \end{pmatrix}$$

לכן  $L(M)$  אינו LTI (לפי הניסוי)



## : addition of a bias (d)

Let  $b \in \mathbb{R}^n$  and  $V \in \mathbb{R}^n$  and  $V_k$  is a vector

$$F(V) = V + b \quad \text{is the mapping}$$

Let  $F(0) = b \neq 0$  then  $F$  is not LTI

Let  $F$  is not LTI (not linear time invariant)

## : Full connected matrix (e)

Let  $V \in \mathbb{R}^n$  and  $V_k$  is a vector

$$M \in \mathbb{R}^{n \times n} \quad F(V) = M \cdot V \quad \text{is the mapping}$$

Let  $V_k$  is a vector and  $V$  is a vector

$$(M V_k)_j = \sum_{t=1}^n M_{jt} \cdot (V_k)_t = \sum_{t=1}^n M_{jt} \cdot V_{t-k \pmod{n}}$$

$$\sum_{t=1}^n M_{j+k \pmod{n}, t} \cdot V_t = (M V)_{j+k \pmod{n}}$$

$$V = \begin{pmatrix} 1 \\ 2 \\ 0 \end{pmatrix}, \quad M = \begin{pmatrix} 0 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 0 \end{pmatrix} \quad \text{is the mapping}$$

$$M \cdot V = \begin{pmatrix} 2 \\ 3 \\ 1 \end{pmatrix}$$

Let

$$V_k = \begin{pmatrix} 0 \\ 1 \\ 2 \end{pmatrix} \quad k=1 \quad \text{is the mapping}$$

$$M \cdot V_k = \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix}$$

$$(M V_k)_2 \neq (M V)_1$$

is the mapping

Let  $F$  is not LTI