



INSTITUTO SUPERIOR TÉCNICO

DEEP LEARNING  
2<sup>nd</sup> QUARTER - 2022/2023

---

## Homework 2

---

*Group 46:*  
Santiago QUINTAS - 93179  
Alina KLINGEL - 105208

*Professor:*  
Mário FIGUEIREDO

15<sup>st</sup> January, 2023

## Declaration of Contribution

We worked on all of the Questions together. In Question 3 we did together in a Discord Call, the other Questions we worked in separate and compared the results. In Question 2 Santiago was leading the process of development, while Alina did so in Question 1.

## Question 1

Annexed file after Q3.

## Question 2

1.

A CNN will have fewer parameters than a fully-connected network with the same input size and the same number of classes, because the output neurons of a CNN will only take into account certain input neurons, while the output neurons of a fully-connected network take into consideration every single input neuron in the layer before. This means that instead of having a set of  $n$  weights for every output neuron, corresponding to  $n$  input neurons, the CNN will have for every output neuron a number of weights lower than  $n$ . This results in a reduction of free parameters. Below there is a visual representation of the density of weights in both networks, [1].

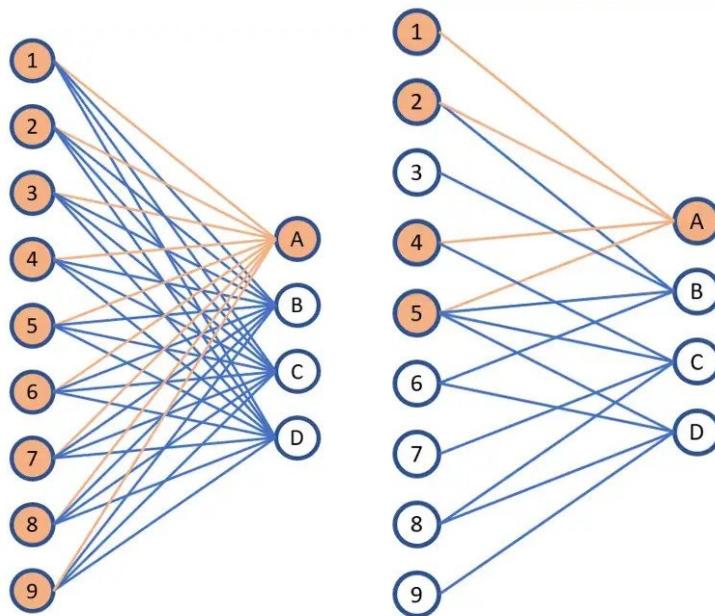


Figure 1: Visual representation of the weight density in a fully connected network (left) and in a CNN (right)

On top of this, a CNN can also share the same weights for different input neurons, so for the example above (Fig 1.), input image of 4x4 pixels, and a filter of 2x2, without weight sharing we use 16 weights, if we apply weight sharing that number is reduced to 4 weights, because the same 4 weights are reused in every iteration of the filter.

## 2.

CNN's make use of the spatial locality of the input pixels, allowing the network to learn patterns, this makes them excel at analysing images and patterns representing letters and numbers, achieving better results than a fully-connected network. While a fully-connected network runs into the problem of overfitting, a CNN does not because of the reduced amount of free parameters, the CNN is also able to achieve better efficiency for big data inputs, such as images, due to the significant decrease in weights used to calculate the output of the neurons. On the opposite side, fully-connected networks lose efficiency when the input data increases by a great amount, since the number of weights that need updating become too big for the network to be efficient.

## 3.

The key benefit of a CNN is the usage of spatial structures to increase the efficiency of learning, when that is taken away, the CNN loses its advantage and the better generalization goes away with it.

In a CNN the spatial structure of an image is kept by making the shape of the pixel vector turn into a 2D array, for example a vector holding 784 pixels becomes a 28x28 array, this is what allows the CNN to detect specific patterns, the filters applied to the image take into account the spatial location of the pixels, and they would lose this ability if they were flattened into a 1 dimension array, without it the CNN wouldn't be able to detect specific patterns in the image and would be the equivalent to a fully-connected network.

## 4.

The default parameters used in the implementation of the CNN were the following:

Number of Epochs	20
Learning Rate	0.01
Batch Size	8
Optimizer	adam

Table 1: Default parameters 1

For different learning rates we achieved different accuracies, shown in the table 2.

The best configuration, according to the accuracy test is with a learning rate of 0.0005

Learning Rate	Final test accuracy
0.01	0.7209
0.0005	0.9581
0.00001	0.8896

Table 2: Accuracy test for different learning rates

Below is the plot of both the training loss and accuracy over the number of epochs of the changes for a learning rate of 0.001, which corresponds to our best configuration.

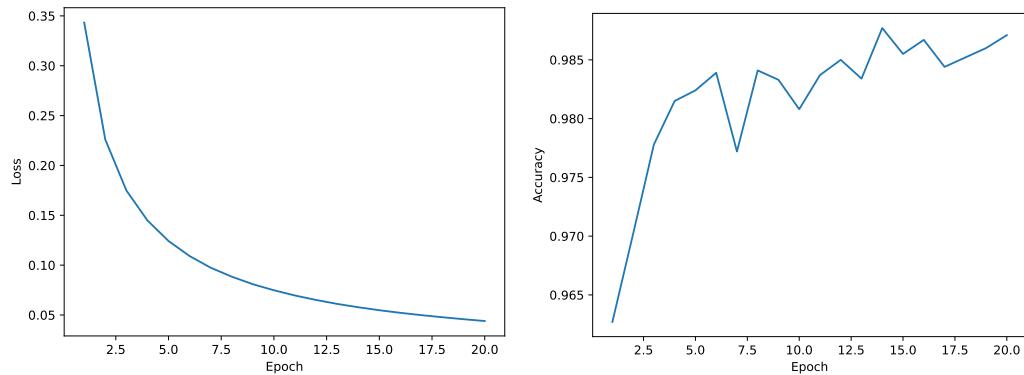


Figure 2: Training loss and accuracy of the CNN for a learning rate of 0.0005

5.

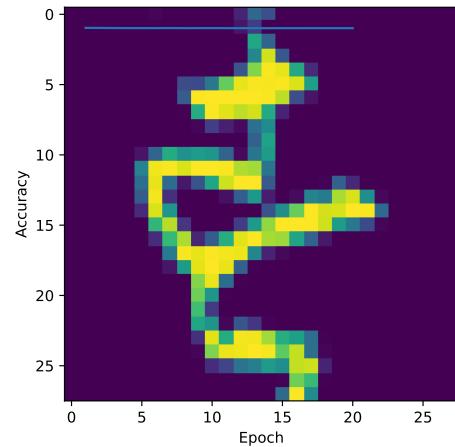


Figure 3: Original training example

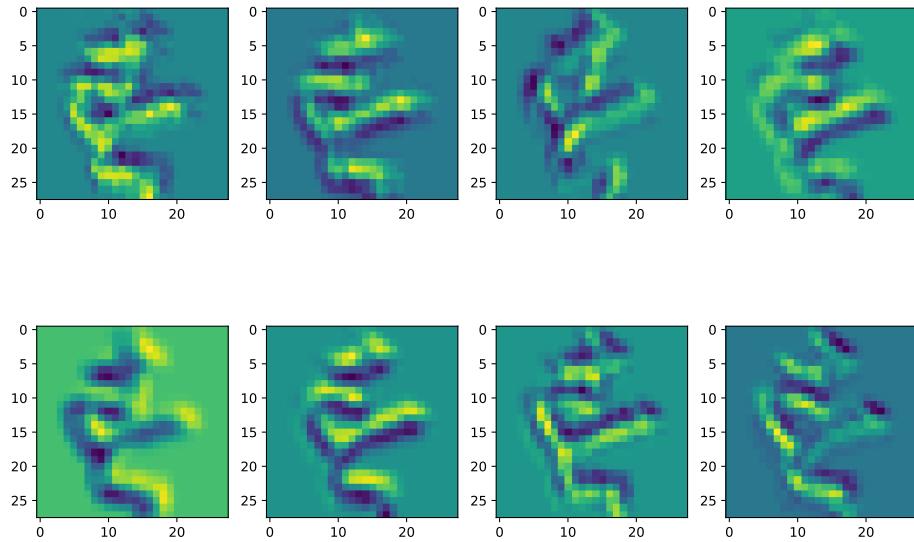


Figure 4: Activation maps of the first convolutional layer for a learning rate of 0.0005

In the activation maps, the highlighted parts correspond to the edges of the original figure, in some maps the bottom horizontal lines are highlighted, in others it's the top horizontal lines, also the vertical lines from both the left and the right of the training example.

## Question 3

1.

a)

The default parameters used in the implementation of machine translation model were the following:

Number of Epochs	50
Learning Rate	0.003
Hidden Size	128
Dropout	0.3
Batch Size	64

Table 3: Default parameters 2

Below is the plot of the validation error over the number of epochs.

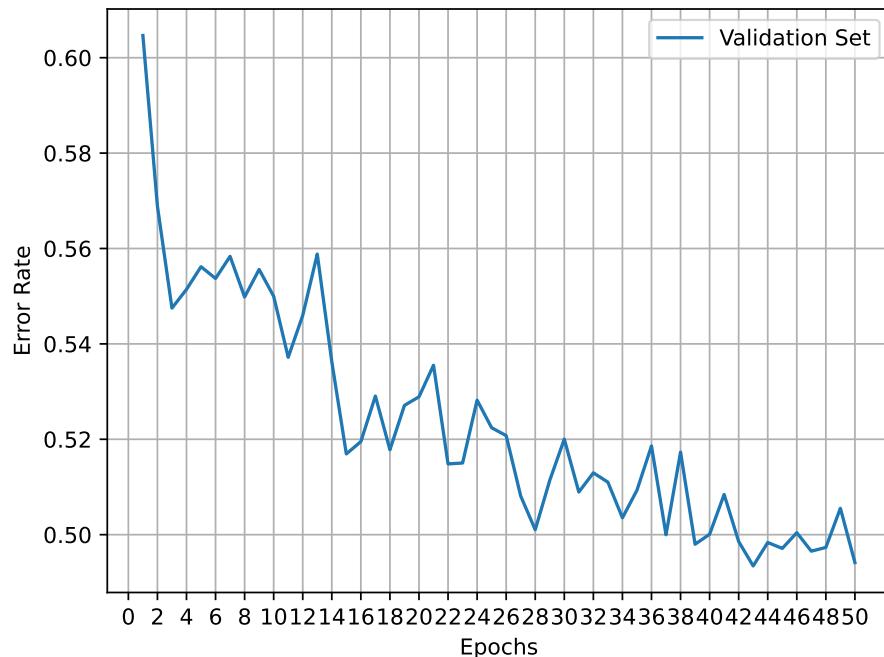


Figure 5: Validation error of the machine translation model as a function of the epoch number

For this configuration we obtained a final error rate of **0.5016**

b)

Below is the plot of the validation error over the number of epochs, with the added attention mechanism and using the same default parameters.

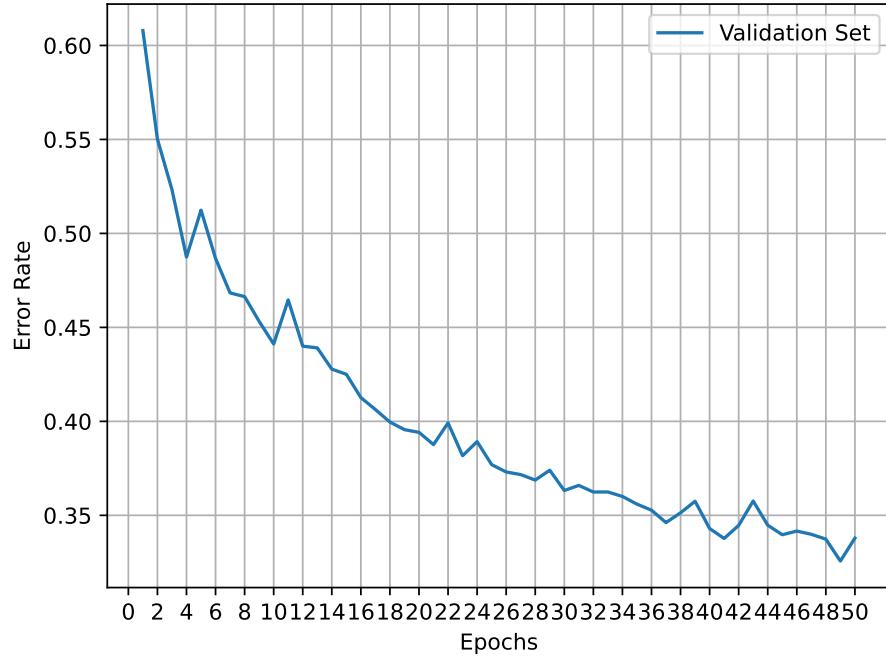


Figure 6: Validation error of the machine translation model, with the added attention mechanism, as a function of the epoch number

For this configuration we obtained a final error rate of **0.3421**

c)

To improve the results, without changing the model architecture, we can use a beam search algorithm in the decoding process as an alternative to the greedy search we are currently using. Instead of keeping a single word, the most probable one, at each time step, the beam algorithm keeps multiple high probability words, this allows for a more flexible decoding process and overall better results.

# Bibliography

- [1] Diego Unzueta, 2021. Deep Learning Fundamentals. Place of publication: Towards Data Science. Available from: <https://towardsdatascience.com/convolutional-layers-vs-fully-connected-layers-364f05ab460b>.

# Question 1 (30 points)

1. (30 points) Consider the convolutional neural network in Fig. 1, used to classify images  $\mathbf{x} \in \mathbb{R}^{H \times W}$  into 3 possible classes. The network has a convolutional layer with a single  $M \times N$  filter, a stride of 1 and no padding, and a ReLU non linearity. This layer is followed by a  $2 \times 2$  max pooling layer with a stride of 2 and no padding, and an output layer with softmax activation.

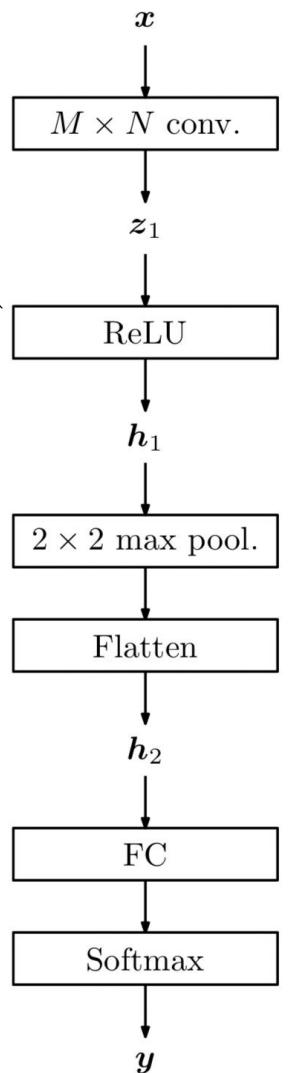
Let  $\mathbf{W}$  denote the filter weights,  $\mathbf{z}$  the result of the convolution (i.e.,  $\mathbf{z} = \text{conv}(\mathbf{W}, \mathbf{x})$ ) and  $\mathbf{h} = \text{ReLU}(\mathbf{z})$ . Let also  $\mathbf{x}' = \text{vec}(\mathbf{x})$  and  $\mathbf{z}' = \text{vec}(\mathbf{z})$  denote the flattened versions of  $\mathbf{x}$  and  $\mathbf{z}$ , respectively.

Given: • image  $\vec{x} \in \mathbb{R}^{H \times W}$  in three possible classes

- 1 conv. layer: • 1  $M \times N$  filter • stride  $S = 1$  • padding  $D = 0$  (no pad),  $(\vec{z} = \text{conv}(\vec{W}, \vec{x}))$
- ReLU = non linear •  $\vec{h} = \text{ReLU}(\vec{z})$

- 1 max pooling layer: •  $2 \times 2$  • stride  $S = 2$  • padding  $D = 0$  (no pad)  $\vec{x}' = \text{vec}(\vec{x})$
- $\vec{z}' = \text{vec}(\vec{z})$

- 1 output layer: softmax activation



a) Dimension of  $\vec{z}$  (= result of convolution)

Idea: Get output (width) of the convolutional layer like in practical 8

$$\text{output width} = \frac{\text{input width} - \text{kernal width} + 2 \times \text{padding width}}{\text{stride}} + 1$$

- input width:  $W$
- kernal width:  $N$
- padding:  $0$
- stride:  $1$

$$\Rightarrow \text{output width} = \frac{W - N + 0}{1} + 1 = W - N + 1$$

• input height:  $H$   
• kernel height:  $M$        $\Rightarrow$  output height =  $\frac{H - M + O}{n} + 1$   
• padding : 0      after convolution  
• stride :  $n$       =  $H - M + 1$

$\Rightarrow$  Dimensions of  $\tilde{x}$  are  $(W - N + 1) \times (H - M + 1)$

- (b) (10 points) Show that there is a matrix  $M \in \mathbb{R}^{H'W' \times HW}$ , with  $H' = H - M + 1$  and  $W' = W - N + 1$ , such that  $\vec{z}' = M\vec{x}$ . 2) What is the general expression for element  $(i, j)$  of  $M$ ?

- matrix  $\vec{M} \in \mathbb{R}^{H'W' \times HW} = \mathbb{R}^{(H-M+1)(W-N+1) \times HW}$

$\downarrow$   
dimensions of  
output of  
convolutional  
layer

$\downarrow$   
dimensions of  
original input  
 $\vec{x}$

$a = \text{Number of rows}$   
 $b = \text{Number of columns}$   
 $a \times b = 5 \times 1$   
 $\Rightarrow 10 \text{ rows} \times \text{columns!}$

Matrix multiplication  
 $n \times p \times m = n \times m$

- vectors:  $\vec{z}'$  and  $\vec{x}'$  = flattened versions of  $\vec{x}$  and  $\vec{z}$   
 $(\text{both } \text{Vec}(\vec{x}))$

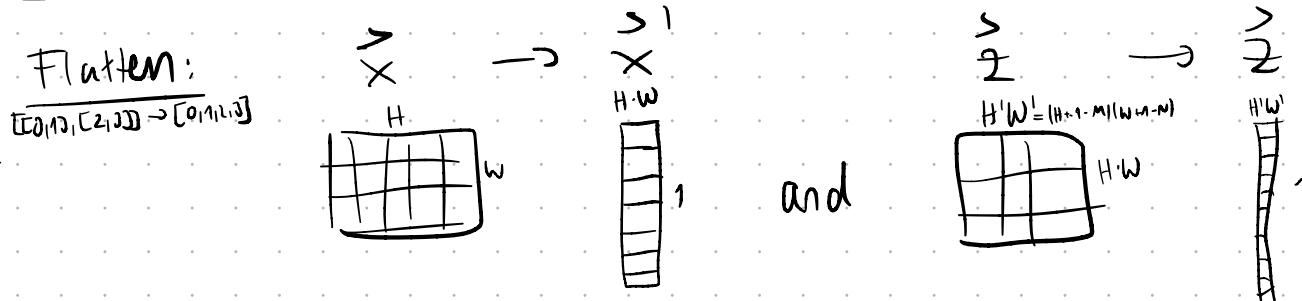
$\uparrow$   
input  
 $\vec{x}' \in \mathbb{R}^{H \times W}$

$\uparrow$   
result of convolution  
 $\mathbb{R}^{(H-M+1) \times (W-N+1)}$

→ go after convolution:

- need to ReLU + maxpool to get to flattened version of  $\vec{z}''$ .
- or can I just flatten  $\vec{x}$  and  $\vec{z}$  to show  $\vec{z}' = M\vec{x}'$

Try to just flatten:



Look at dimensions:

$\vec{x}' : \overset{\text{rows}}{H \cdot W} \times \overset{\text{columns}}{1}$   $\Rightarrow$  with  $\vec{z}' = \vec{M} \cdot \vec{x}'$  for the product to be possible  $\vec{M}$  needs the dimension

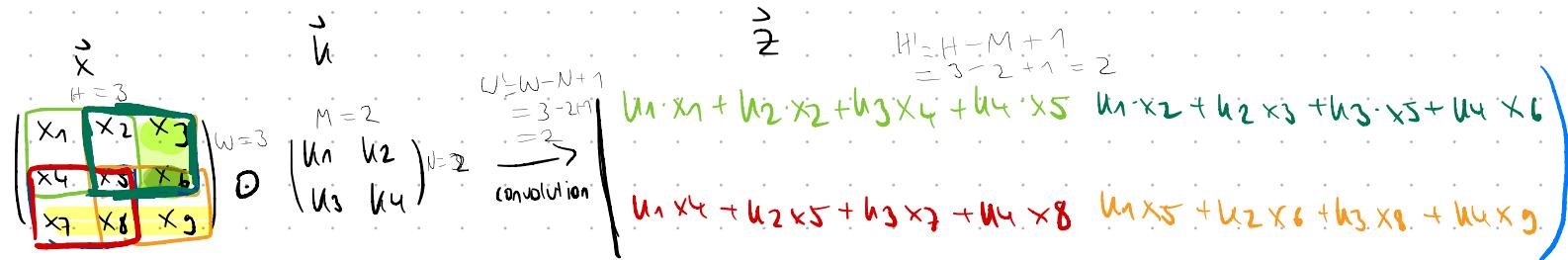
of  $H'W' \times HW$   $[H'W' \times 1 = (H'W') \times (HW)]$

$\Rightarrow$  so we showed, that the matrix to be multiplied with  $\vec{x}'$  needs the dimensions  $H'W' \times HW$  in order to get the vector of  $\vec{z}'$  with dimension  $H'W' \times 1$

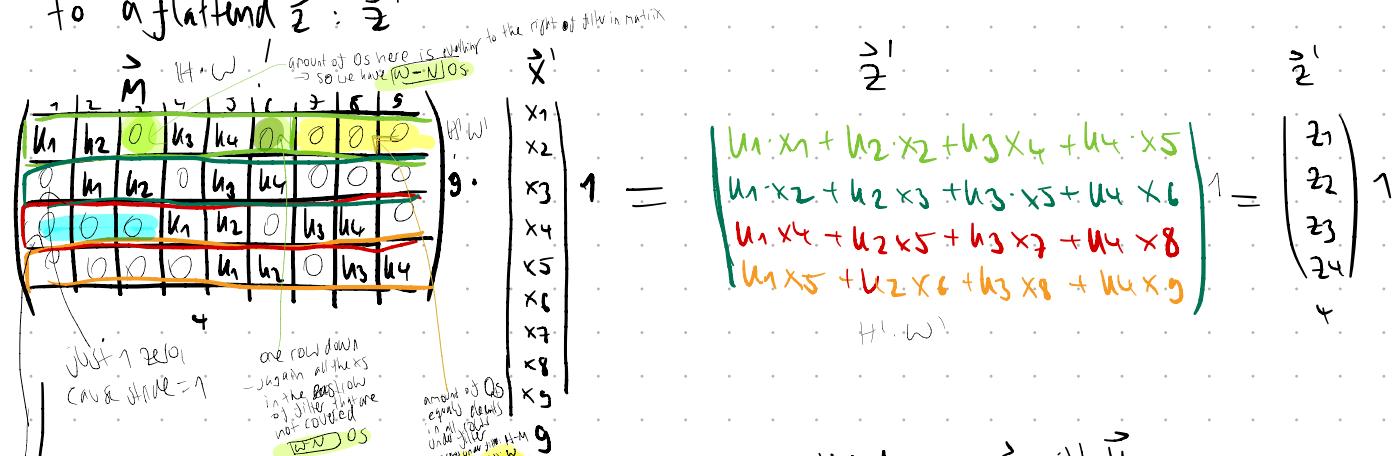
2) Find General expression for a element  $(i,j)$  of  $\hat{M}$

2.1 Look at a small example: What is happening?

- Convolution  $\hat{z}$  for a 2d input  $\hat{x}$  and 2D kernel  $\hat{v}$  (with stride 1)  
 $\Rightarrow$  Sliding window  $\hat{v}$  over  $\hat{x}$  to get matrix  $\hat{z}$



- It is possible to construct a Matrix  $\tilde{M}$  that leads as product with the flattened  $\tilde{x}$ :  $\tilde{x}^T \tilde{M} \tilde{x}$  to a flattened  $\tilde{z} : \tilde{z}^T$



so this gives same result as above where we slided over  $\vec{x}$  with  $\vec{w}$   
 we shifted one down  $\rightarrow$  so one full offset of  $\vec{w}$ :  $\vec{w} \vec{w}$

22. Now try to generalize this (for stride 1!)

Input: HxW - Matrix

Filter (Kernel):  $M \times N$ -Matrix

convolutet Matrix:  $H' \times W'$ -Matrix

$$\begin{pmatrix} x_{11} & \cdots & x_{1N} \\ \vdots & \ddots & \vdots \\ x_{H1} & \cdots & x_{HN} \end{pmatrix} \odot \begin{pmatrix} k_{11} & \cdots & k_{1N} \\ \vdots & \ddots & \vdots \\ k_{H1} & \cdots & k_{HN} \end{pmatrix} \rightarrow$$

$\vec{x} \in \mathbb{R}^{H \times W}$        $\vec{k} \in \mathbb{R}^{M \times N}$

mech → let's just do M like this (no brain dies.)

$\vec{z}e$   $H' \times W'$

↳ now as  $\vec{m} \cdot \vec{x}' = \frac{1}{2}$

## Examples



- one row of  $\{g_i\}$  with elements  $g_{ij}$  with  $1 \leq j \leq M \cdot n$

$$g_j = \begin{cases} k, & \text{if } k \in N \\ 0, & \text{otherwise} \end{cases}$$

$\Rightarrow$   $k = (j-1) \bmod (N-1)$

$b = ((j-1) \bmod (N-1)) + 1$

$\Rightarrow j \bmod (N-1) + 1$

$$\begin{array}{l} \text{91. } a = ((-1)^2) \operatorname{div} 11 + 1 = 1 \\ b = (-1)^{11} \operatorname{mod} 11 - 1 = 1 \end{array} \quad \begin{array}{l} \text{92. } a = (2-1) \operatorname{div} 3 + 1 = 1 \\ b = (2-1) \operatorname{mod} 3 + 1 = 2 \end{array} \quad \begin{array}{l} \text{93. } a = (3-1) \operatorname{div} 3 - 1 = 1 \\ b = (3-1) \operatorname{mod} 3 + 1 = 4 \end{array}$$

Shift :=  $i \bmod H^1$        $H^1 = H - M + 1 = 3 - 2 + 1 = 2$        $i = \text{rows}$   
 $i=1: \text{shift} = 1 \bmod 2 = 1$        $i=4: \text{shift} = 4 \bmod 2 = 0$        $M_{H^1, i=1, j=1}$   
 $i=2: \text{shift} = 2 \bmod 2 = 0$        $i=5: \text{shift} = 5 \bmod 2 = 1$        $M_{H^1, i=2, j=0}$   
 $i=3: \text{shift} = 3 \bmod 2 = 1$        $i=6: \text{shift} = 6 \bmod 2 = 0$        $M_{H^1, i=3, j=1}$

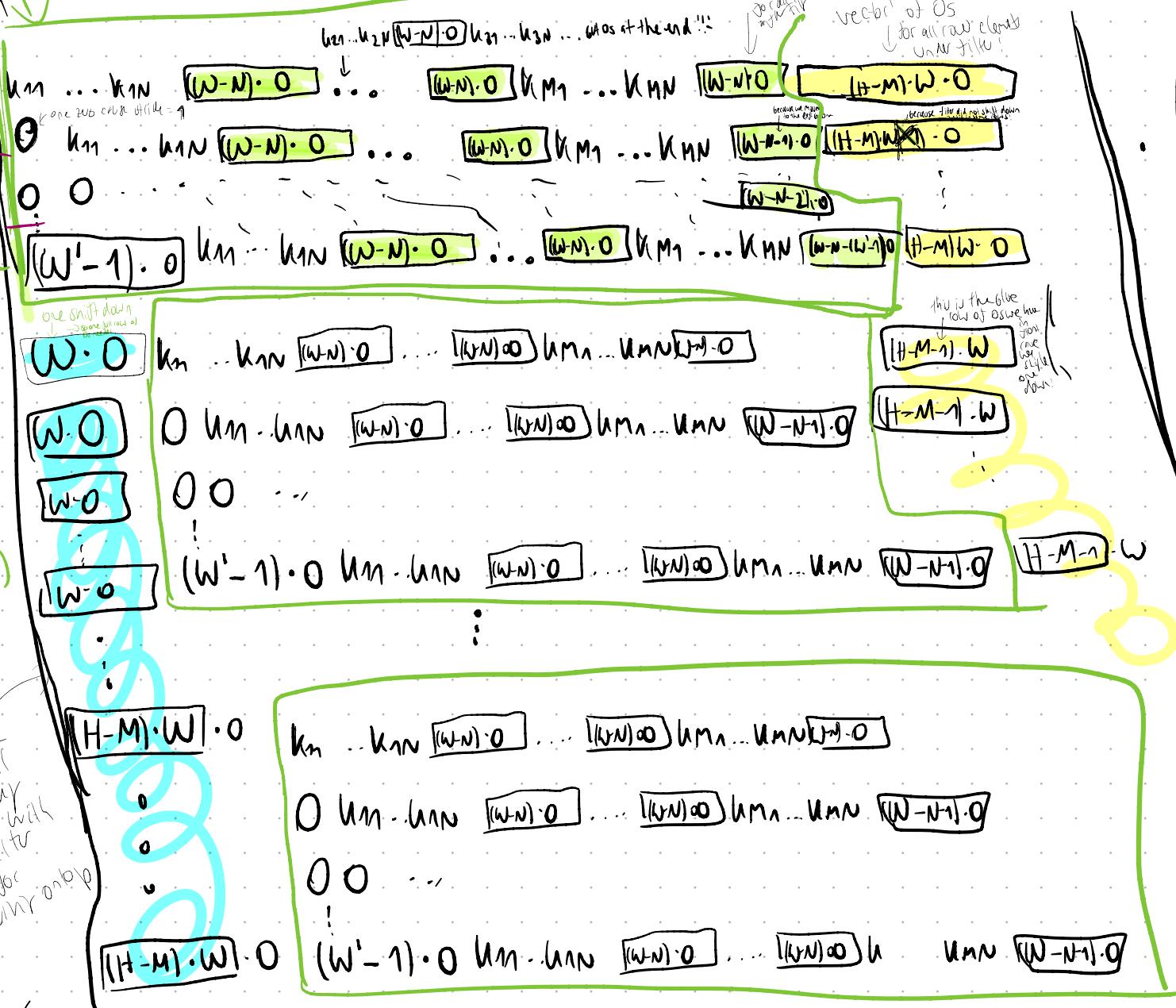
in this box:  
result of shifting  $W-1$  times filter to the right  
(until filter hits right end of  $\vec{x}$  matrix)

↓ cause we don't want to count it twice

$Gaps = 0$   
shift  
 $W-1$  times to  
right  
(or  $W-n$   
and without  
position).

$M :=$   
 $Gaps = 1$   
(repeat  
this  $H-1$   
times  
(so shift it one  
position down  
and same  
thing again))

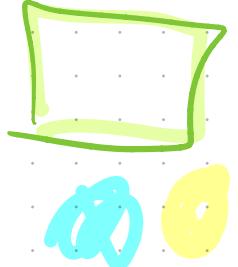
now we  
are all  
the way  
down with  
the filter  
→ 0s for  
every row



$$\vec{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_{1W} \\ x_{21} \\ x_{22} \\ \vdots \\ x_{H1} \\ x_{HW} \end{pmatrix}$$

rows ↓ columns ↓

$$M_{ij} = \begin{cases} K_{ab} & \text{if } i = \frac{a}{b} \text{ with } a \in \\ & j \in \end{cases} ? ?$$



• Box  $\boxed{\square}$  with just  $M_{ij}$ 's and 0's:

$\downarrow$  row column  
 $\downarrow j$

$\therefore$  matrix  $G_1$  with elements  $g_{ij}$

$$a \bmod b = \text{Rest}\left(\frac{a}{b}\right)$$

$$1 \bmod 3$$

$$1:3=0 \text{ Rest } 1$$

• Dimensions of Matrix  $G_1$  :

$$\text{number of rows: } W^1 = W - N + 1$$

$$\rightarrow 1 \leq i \leq W^1$$

$$\text{number of columns: } M \cdot N + M \cdot (W-N) = M(N+W-N) = M \cdot W$$

$W^1$

$M \cdot W \quad W^1 \times M \cdot W$

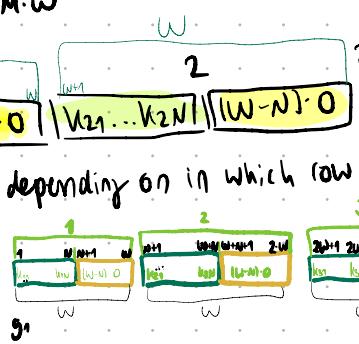
• Each row:

$\times$  sequence of  $\boxed{k_{11} \dots k_{1N}} \boxed{(W-N) \cdot 0} \boxed{k_{21} \dots k_{2N}} \boxed{(W-N) \cdot 0} \dots \boxed{k_{M1} \dots k_{MN}} \boxed{(W-N) \cdot 0}$

(+ 0's to left and right of box depending on in which row we are)

\*  $k$ -sequences: always  $N$  long

0-sequence: always  $(W-N)$  long



\* One row of  $G$  with elements  $g_{ij}$  with  $1 \leq j \leq M \cdot W$

the elements  $g_{ij}$  at this row can be defined as:

$$g_{ij} = \begin{cases} K_{ab}, \text{ if } b \leq N & a := (j-1) \div W + 1 \\ 0 & \text{otherwise} \end{cases}$$



ONE ROW OF  $G$  WITH ELEMENTS  $g_{ij}$  TO MATRIX

1 --- N

g1 g2 g3 g4 g5 g6

use this concept to find out how many zeros there are in a row of our given matrix  $G$ .  
Imagine the square elements as the value in a matrix.

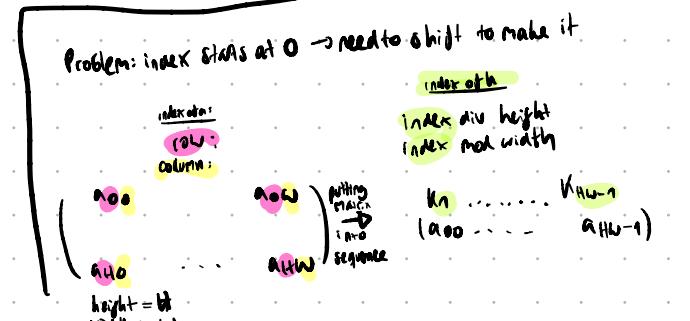
• Now figure out the 0's surrounding  $\boxed{\square}$

\* Behind every row of the first  $\boxed{\square}$  there are  $(H-M) \cdot W$  0's

\* For the second  $\boxed{\square}$  before each row there are  $(W-1)W$  0's and behind  $(H-N-1)W$  0's

\* For the third  $\boxed{\square}$  before each row there are  $2W$  0's and behind  $(H-m-2)W$  0's

How to segmentize matrices



\* Finally for the last  $\square$  in  $M$  before each row  
 there are  $(H-W) \cdot W$  0s and Behind are none

- In each  $\square$  are  $W$  positions / shifts

Every row of  $\square$  and therefore every row of  $M$   
 corresponds to one unique position of the filter kernel

↳ define each unique position as g-pos

$$\boxed{g\text{-pos} := i \bmod H^1} \quad i = 0, \dots, H^1$$

↓  
 0-blocks to the left

$$\Rightarrow M(i, j) = \begin{cases} 0 & \text{if } j \leq g\text{-pos} \cdot W \\ \text{or } j > M \cdot W + g\text{-pos} \cdot W \\ g_j & \text{otherwise} \end{cases}$$

↑  
 0-block to the right

$$M \cdot W + g\text{-pos} \cdot W$$

size of gram block      size of the 0-block

$$W\text{-pos} := j - g\text{-pos} \cdot W$$

shift it to the very left to be able to  
 get column index of  
 the last element!

(c) (10 points) Indicate the number of parameters in the network. Compare that number with the number of parameters in a network where the convolutional and max pooling layers are replaced by a fully connected layer yielding an output with the same dimension as  $h_2$ . Ignore the bias terms.

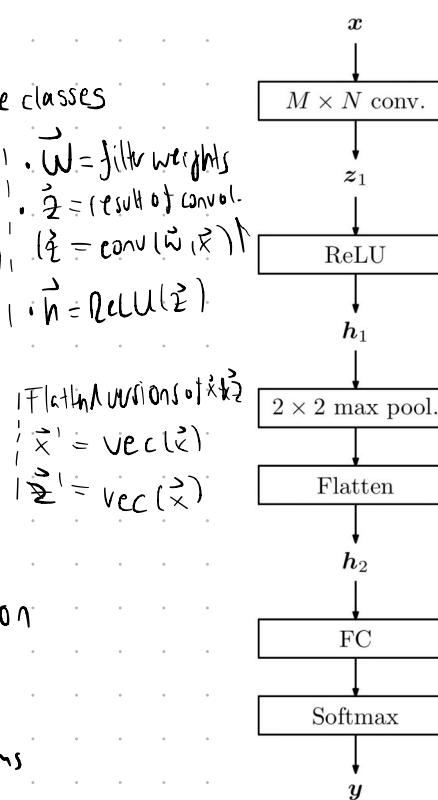
Probably similar to practical 8, Question 2.3 and 2.4

→ Do they want all parameters or just the trainable ones? I will just do the trainable ones

## 1 CNN-Network

Given: image  $\vec{x} \in \mathbb{R}^{H \times W}$  in three possible classes

- 1 conv. layer:  $1 \cdot M \times N$  filter  
· stride  $S = 1$   
· padding  $D = 0$  (no pad)  
· ReLU = non linear
- 1 MAX pooling layer:  $2 \times 2$   
· stride  $S = 2$   
· padding  $D = 0$  (no pad)
- 1 output layer: softmax activation



### Number of parameters

input layer: nothing to learn, just provide input image's shape

# params  
0

conv. layer: CNN learns here → have weight

matrices, number of trainable weights:  
 $\# \text{of filters} \times (\text{filter width} \times \text{filter height} \times \text{# of channels}) + b^{\text{bias}}$

(from practical 8.2.1)

$1 \times ((M \times N \times 1) + 0)$

↑  
assume no bias?

pool layer: no learnable params, because we just calculate a specific number → no backprop learning

0

↑ previous

$\dim(h_2) \cdot 3$

↑ 3 outputs = dim(y)

$$= \frac{H' \cdot W'}{4} \cdot 3$$

FC layer: fully connected layer

→ has learnable params  
→ highest number: because every neuron is connected with every neuron:

# of neurons at previous layer

\* # of neurons of current layer

(connects all neurons from output of max pooling to the final 3 neurons of the classification layer)  
↳ 3 classes

calculation of  $h_2$ :

• convolution:

→ dim of  $z_1$ :  $H' \times W' = (H - M + 1) \times (W - N + 1)$  (see 2.1.1)

• ReLU changes nothing on dim (because output of 0)

→ dim of  $h_1$  = dim of  $z_1$  =  $H' \times W'$

•  $2 \times 2$  Max pool:  $\frac{H - M + 1}{2} \times \frac{W - N + 1}{2}$

flattening

→ dim of  $h_2$ :  $\frac{H - M + 1}{2} \cdot \frac{W - N + 1}{2} = \frac{(H - M + 1)(W - N + 1)}{4} = \frac{H' \cdot W'}{4}$

$$\Rightarrow \text{Total # of trainable params} | N_{\text{CNN}} = M \cdot N + \frac{H' \cdot W'}{4} \cdot 3 |$$

## ② Fully connected layers : output of dim (h2)!

# params

FC-layer 1 :-  $\dim(\text{input}) \cdot \dim(h_2)$

$$= (H \cdot W) \cdot \frac{H \cdot W}{4}$$

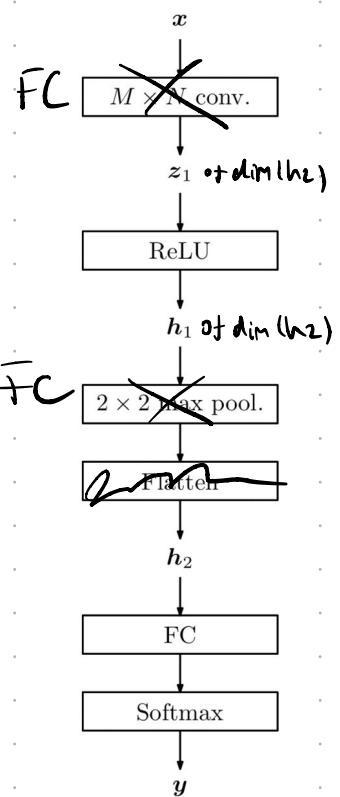
FC-layer 2 :-  $\dim(h_2) \cdot \dim(\text{class})$   
or  $y$

$$= \frac{H \cdot W}{4} \cdot 3$$

=> total # of trainable param.

$$\boxed{N_{FC} = (H \cdot W) \cdot \frac{H \cdot W}{4} + \frac{H \cdot W}{4} \cdot 3}$$

$$= \frac{H \cdot W}{4} (H \cdot W + 3)$$



## ③ Comparison

$$N_{CNN} = M \cdot N + \frac{H \cdot W}{4} \cdot 3$$

$$N_{FC} = H \cdot W \frac{H \cdot W}{4} + \frac{H \cdot W}{4} \cdot 3$$

1)  $\frac{H \cdot W}{4} > 1$  ← we assume this (when a  $\frac{1}{4}$  of  $n \times n$  is weird!)

2)  $MN < H \cdot W$  (because filter size is smaller than image for sure!)  $(M \cdot N) < (H \cdot W)$

$$\Rightarrow MN < H \cdot W \cdot \frac{H \cdot W}{4}$$

$$\Rightarrow N_{CNN} < N_{FC}$$

2. (5 points) Assume now that, instead of using convolutions, the same  $H \times W$  image  $x$  is flattened into a sequence  $x' = \text{vec}(x)$  of length  $HW$  and goes through a single-head self-attention layer with  $1 \times 1$  projection matrices  $W_Q = W_K = W_V = 1$ , without any positional encodings. Write the expression for the attention probabilities and attention output as a function of  $x'$ .

Same as in Practical 10, Question 1.2  $\rightarrow$  just a lot easier  
+ lecture

query:  $Q = \underset{\substack{\uparrow \\ \text{input matrix}}}{X} \cdot W_Q = X' \cdot W_Q = X'$

key:  $U = \underset{\substack{\uparrow \\ \text{input matrix}}}{X} \cdot W_K = X' \cdot 1 = X'$

value:  $V = \underset{\substack{\uparrow \\ \text{input matrix}}}{X} \cdot W_V = X' \cdot 1 = X'$

(matrices)

attention probabilities:  $p = \text{softmax} \left( \frac{Q U^T}{\sqrt{d_k}} \right)$

$$P = \text{softmax} \left( \frac{X' X'^T}{\sqrt{d_k}} \right)$$

$\overbrace{\text{softmax}}^{\text{sum}(u_i) = 1}$   
 $\overbrace{\text{problem:}}^{\text{single-head}} \quad \begin{array}{l} \text{Attention} \\ = \text{scaled dot product} \\ \text{attention?} \end{array}$

attention:  $Z = P \cdot V$

$$= \text{softmax} \left( \frac{X' X'^T}{\sqrt{d_k}} \right) \cdot X'$$

Scaled Dot-Product Attention  
 Scaling before doing softmax  
 to avoid overflow

**Problem:** As  $d_k$  gets large, the variance of  $q^T k$  increases, the softmax gets very peaked, hence its gradient gets smaller.

**Solution:** scale by length of query/key vectors:

$$Z = \text{softmax} \left( \frac{Q K^T}{\sqrt{d_k}} \right) V.$$

*so softmax would not be so peaked*  
*doesn't matter if it's linear?*