

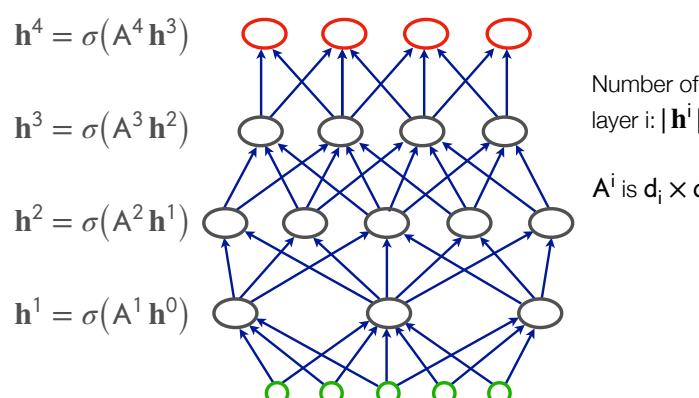
COS324: INTRODUCTION TO MACHINE LEARNING

Prof. Yoram Singer



Topic: Convolutional Neural Networks

© 2020 YORAM SINGER



© 2020 YORAM SINGER 2

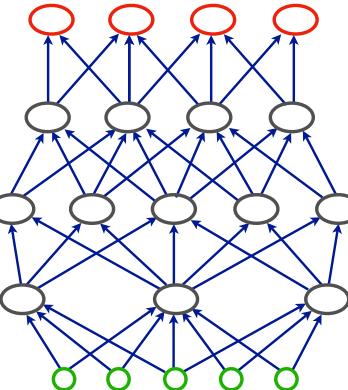
(Almost) Fully Connected Net

$$\mathbf{h}^4 = \sigma(\mathbf{A}^4 \mathbf{h}^3)$$

$$\mathbf{h}^3 = \sigma(\mathbf{A}^3 \mathbf{h}^2)$$

$$\mathbf{h}^2 = \sigma(\mathbf{A}^2 \mathbf{h}^1)$$

$$\mathbf{h}^1 = \sigma(\mathbf{A}^1 \mathbf{h}^0)$$



© 2020 YORAM SINGER 2

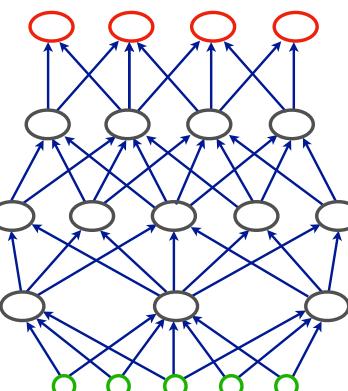
(Almost) Fully Connected Net

$$\mathbf{h}^4 = \sigma(\mathbf{A}^4 \mathbf{h}^3)$$

$$\mathbf{h}^3 = \sigma(\mathbf{A}^3 \mathbf{h}^2)$$

$$\mathbf{h}^2 = \sigma(\mathbf{A}^2 \mathbf{h}^1)$$

$$\mathbf{h}^1 = \sigma(\mathbf{A}^1 \mathbf{h}^0)$$



Number of neurons in layer i : $|\mathbf{h}^i| = d_i$

\mathbf{A}^i is $d_i \times d_{i-1}$ matrix

Need to constrain each \mathbf{A}^i to save space & runtime

© 2020 YORAM SINGER 2

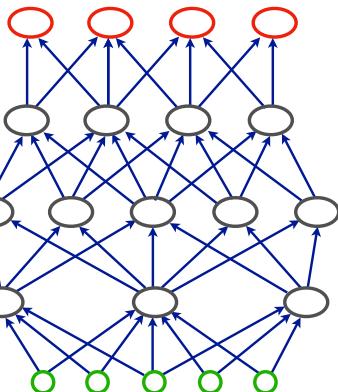
(Almost) Fully Connected Net

$$\mathbf{h}^4 = \sigma(\mathbf{A}^4 \mathbf{h}^3)$$

$$\mathbf{h}^3 = \sigma(\mathbf{A}^3 \mathbf{h}^2)$$

$$\mathbf{h}^2 = \sigma(\mathbf{A}^2 \mathbf{h}^1)$$

$$\mathbf{h}^1 = \sigma(\mathbf{A}^1 \mathbf{h}^0)$$



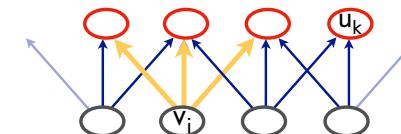
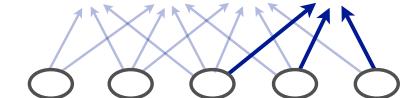
Number of neurons in layer i : $|\mathbf{h}^i| = d_i$

\mathbf{A}^i is $d_i \times d_{i-1}$ matrix

Need to constrain each \mathbf{A}^i to save space & runtime

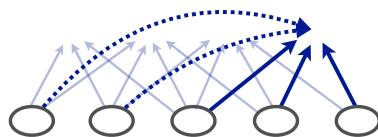
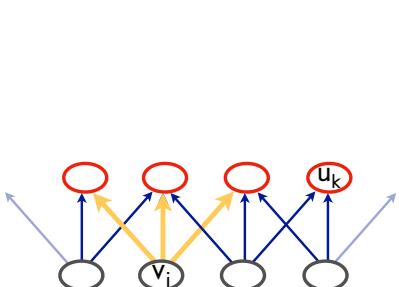
© 2020 YORAM SINGER 2

Sparsity Constraints

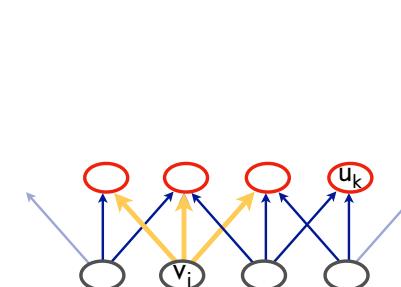


© 2020 YORAM SINGER 3

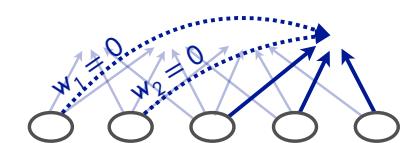
Sparsity Constraints



Sparsity Constraints



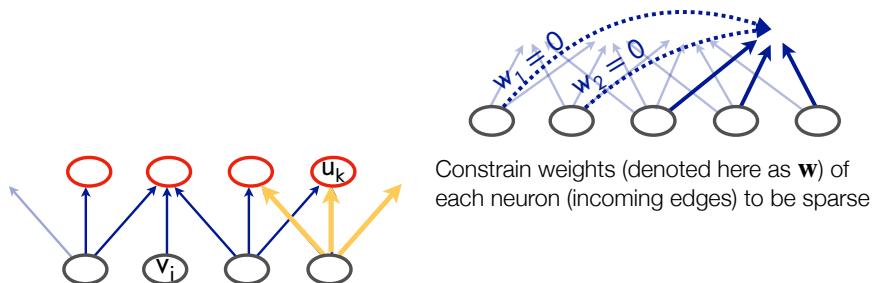
Constrain weights (denoted here as \mathbf{w}) of each neuron (incoming edges) to be sparse



© 2020 YORAM SINGER 3

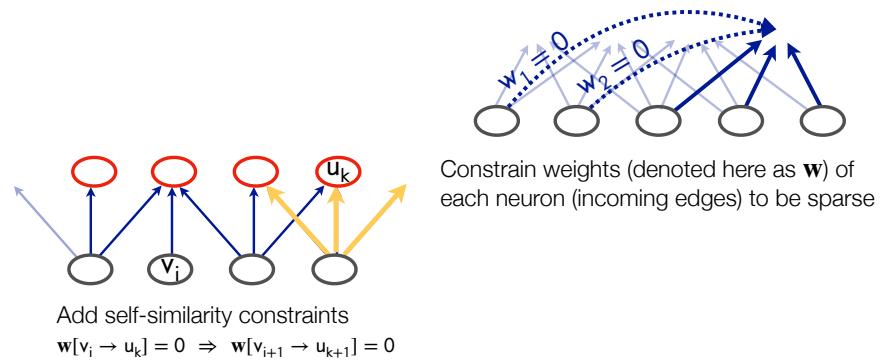
© 2020 YORAM SINGER 3

Sparsity Constraints



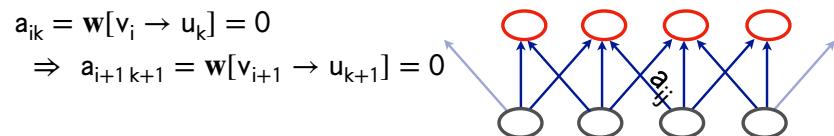
© 2020 YORAM SINGER 3

Sparsity Constraints



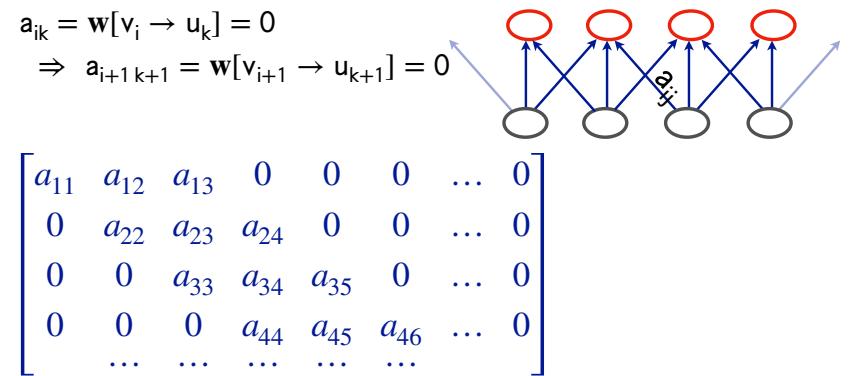
© 2020 YORAM SINGER 3

Structural Self-Similarity



© 2020 YORAM SINGER 4

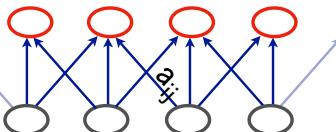
Structural Self-Similarity



© 2020 YORAM SINGER 4

Structural Self-Similarity

$$a_{ik} = w[v_i \rightarrow u_k] = 0 \\ \Rightarrow a_{i+1, k+1} = w[v_{i+1} \rightarrow u_{k+1}] = 0$$

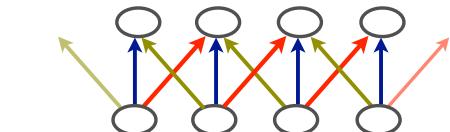


$$\begin{bmatrix} a_{11} & a_{12} & a_{13} & 0 & 0 & 0 & \dots & 0 \\ 0 & a_{22} & a_{23} & a_{24} & 0 & 0 & \dots & 0 \\ 0 & 0 & a_{33} & a_{34} & a_{35} & 0 & \dots & 0 \\ 0 & 0 & 0 & a_{44} & a_{45} & a_{46} & \dots & 0 \\ \dots & 0 \end{bmatrix}$$

Can we further save?

© 2020 YORAM SINGER 4

Parameter Tying

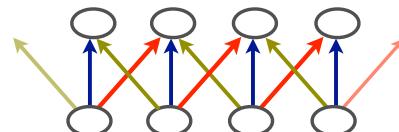


$$\begin{bmatrix} a_{11} & a_{12} & a_{13} & 0 & 0 & 0 & \dots & 0 \\ 0 & a_{22} & a_{23} & a_{24} & 0 & 0 & \dots & 0 \\ 0 & 0 & a_{33} & a_{34} & a_{35} & 0 & \dots & 0 \\ 0 & 0 & 0 & a_{44} & a_{45} & a_{46} & \dots & 0 \\ \dots & 0 \end{bmatrix}$$

© 2020 YORAM SINGER 5

Parameter Tying

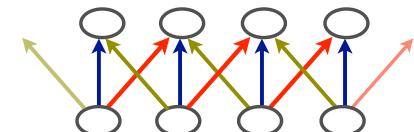
$$\begin{bmatrix} a_{11} & a_{12} & a_{13} & 0 & 0 & 0 & \dots & 0 \\ 0 & a_{22} & a_{23} & a_{24} & 0 & 0 & \dots & 0 \\ 0 & 0 & a_{33} & a_{34} & a_{35} & 0 & \dots & 0 \\ 0 & 0 & 0 & a_{44} & a_{45} & a_{46} & \dots & 0 \\ \dots & 0 \end{bmatrix}$$



© 2020 YORAM SINGER 5

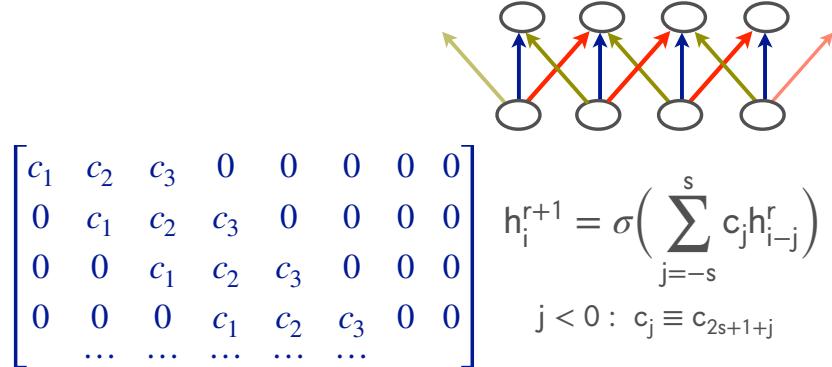
Parameter Tying

$$\begin{bmatrix} c_1 & c_2 & c_3 & 0 & 0 & 0 & 0 & 0 \\ 0 & c_1 & c_2 & c_3 & 0 & 0 & 0 & 0 \\ 0 & 0 & c_1 & c_2 & c_3 & 0 & 0 & 0 \\ 0 & 0 & 0 & c_1 & c_2 & c_3 & 0 & 0 \\ \dots & 0 \end{bmatrix}$$



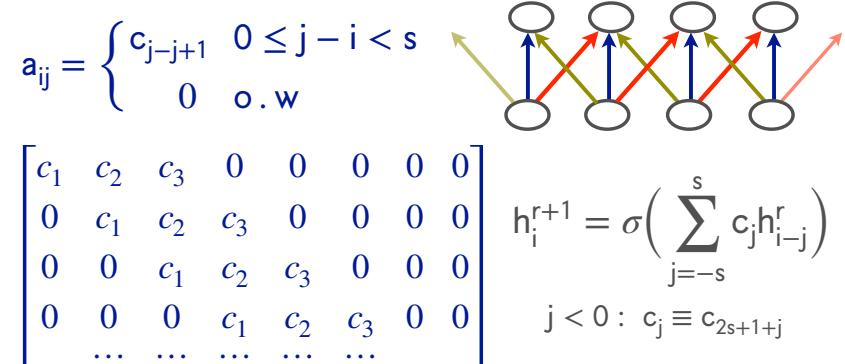
© 2020 YORAM SINGER 6

Parameter Tying



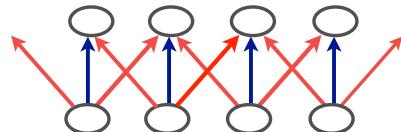
© 2020 YORAM SINGER 6

Parameter Tying



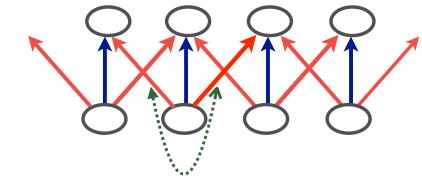
© 2020 YORAM SINGER 6

Parameter Symmetrization



© 2020 YORAM SINGER 7

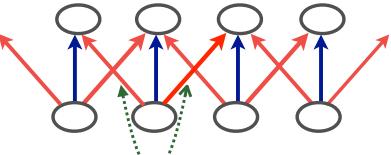
Parameter Symmetrization



© 2020 YORAM SINGER 7

Parameter Symmetrization

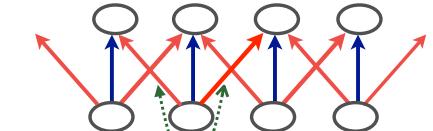
$$\begin{bmatrix} c_1 & c_2 & c_1 & 0 & 0 & 0 & 0 & 0 \\ 0 & c_1 & c_2 & c_1 & 0 & 0 & 0 & 0 \\ 0 & 0 & c_1 & c_2 & c_1 & 0 & 0 & 0 \\ 0 & 0 & 0 & c_1 & c_2 & c_1 & 0 & 0 \\ \dots & \dots \end{bmatrix}$$



Parameter Symmetrization

$$a_{ij} = \begin{cases} c_{|i-j|} & |j - i| < \frac{s}{2} \\ 0 & \text{o.w} \end{cases}$$

$$\begin{bmatrix} c_1 & c_2 & c_1 & 0 & 0 & 0 & 0 & 0 \\ 0 & c_1 & c_2 & c_1 & 0 & 0 & 0 & 0 \\ 0 & 0 & c_1 & c_2 & c_1 & 0 & 0 & 0 \\ 0 & 0 & 0 & c_1 & c_2 & c_1 & 0 & 0 \\ \dots & \dots \end{bmatrix}$$

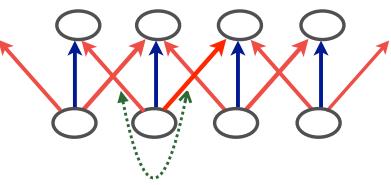


© 2020 YORAM SINGER 7

Parameter Symmetrization

$$a_{ij} = \begin{cases} c_{|i-j|} & |j - i| < \frac{s}{2} \\ 0 & \text{o.w} \end{cases}$$

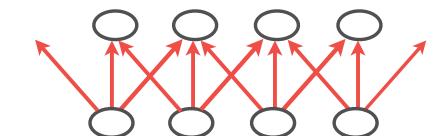
$$\begin{bmatrix} c_1 & c_2 & c_1 & 0 & 0 & 0 & 0 & 0 \\ 0 & c_1 & c_2 & c_1 & 0 & 0 & 0 & 0 \\ 0 & 0 & c_1 & c_2 & c_1 & 0 & 0 & 0 \\ 0 & 0 & 0 & c_1 & c_2 & c_1 & 0 & 0 \\ \dots & \dots \end{bmatrix}$$



$$h_i^{r+1} = \sigma\left(\sum_{j=0}^s c_j \frac{h_{i-j}^r + h_{i+j}^r}{2}\right)$$

© 2020 YORAM SINGER 7

Averaging Neuron

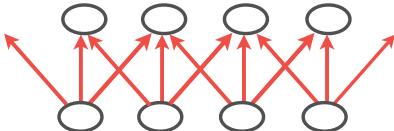


© 2020 YORAM SINGER 8

Averaging Neuron

$$a_{ij} = \begin{cases} \frac{1}{2s+1} & |j-i| < \frac{s}{2} \\ 0 & \text{o.w} \end{cases}$$

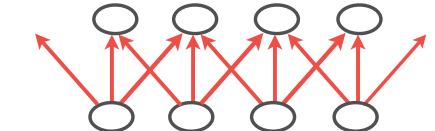
$$\begin{bmatrix} c & c & c & 0 & 0 & 0 & 0 & 0 \\ 0 & c & c & c & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & c & c & c & 0 & 0 \\ 0 & 0 & 0 & 0 & c & c & c & 0 \\ 0 & 0 & 0 & 0 & 0 & c & c & c \end{bmatrix}$$



Averaging Neuron

$$a_{ij} = \begin{cases} \frac{1}{2s+1} & |j-i| < \frac{s}{2} \\ 0 & \text{o.w} \end{cases}$$

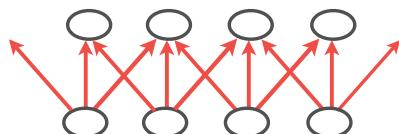
$$\begin{bmatrix} c & c & c & 0 & 0 & 0 & 0 & 0 \\ 0 & c & c & c & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & c & c & c & 0 & 0 \\ 0 & 0 & 0 & 0 & c & c & c & 0 \\ 0 & 0 & 0 & 0 & 0 & c & c & c \end{bmatrix}$$



$$h_i^{r+1} = \sigma\left(\frac{1}{2s+1} \sum_{j=-s}^s h_{i-j}^r\right)$$

© 2020 YORAM SINGER 8

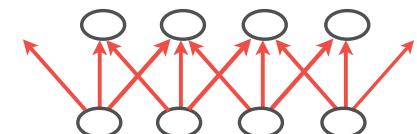
Max-Pooling (Neuron)



© 2020 YORAM SINGER 9

Max-Pooling (Neuron)

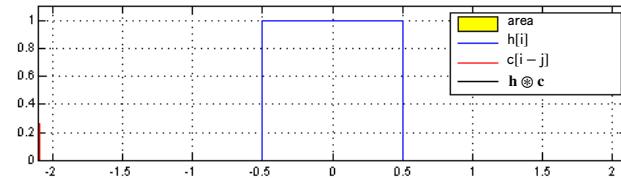
$$h_i^{r+1} = \sigma\left(c \max_{j=-s}^s h_{i-j}^r\right)$$



© 2020 YORAM SINGER 9

Convolution

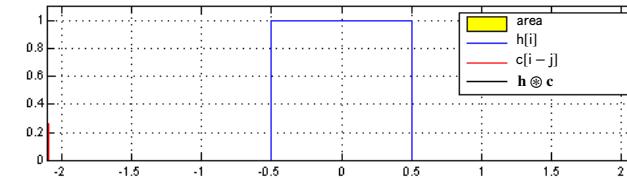
$$h_i^{r+1} = \sigma \left(\underbrace{\sum_{j=-s}^s c_j h_{i-j}^r}_{z_i} \right) \quad (\mathbf{c} \circledast \mathbf{h})_i = \sum_{j=-s}^s c_j h_{i-j}$$



© 2020 YORAM SINGER 10

Convolution

$$h_i^{r+1} = \sigma \left(\underbrace{\sum_{j=-s}^s c_j h_{i-j}^r}_{z_i} \right) \quad (\mathbf{c} \circledast \mathbf{h})_i = \sum_{j=-s}^s c_j h_{i-j}$$



© 2020 YORAM SINGER 10

Boundary Conditions

For $i \leq s$ and $i > r - s$:

Convolution is not uniquely defined

Boundary Conditions

For $i \leq s$ and $i > r - s$:

Convolution is not uniquely defined

h

0.1	7	-1	2	3
-----	---	----	---	---

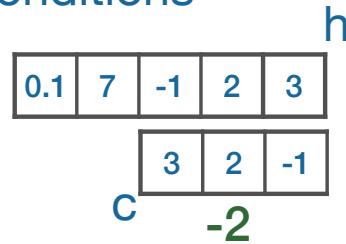
© 2020 YORAM SINGER 11

© 2020 YORAM SINGER 11

Boundary Conditions

For $i \leq s$ and $i > r - s$:

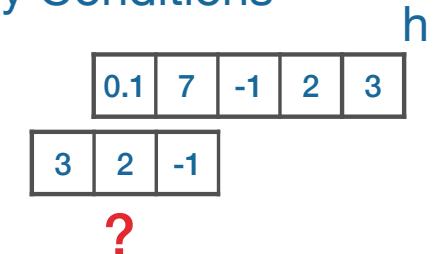
Convolution is not uniquely defined



Boundary Conditions

For $i \leq s$ and $i > r - s$:

Convolution is not uniquely defined



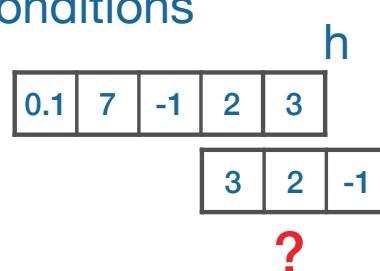
© 2020 YORAM SINGER 11

© 2020 YORAM SINGER 11

Boundary Conditions

For $i \leq s$ and $i > r - s$:

Convolution is not uniquely defined



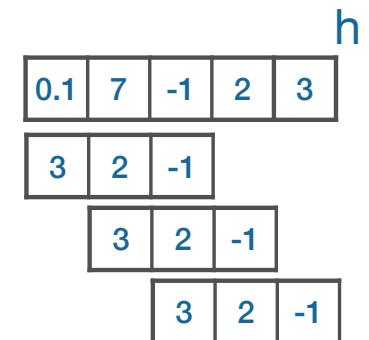
Valid Only Convolution

For $i \leq s$ and $i > r - s$:

Convolution is not uniquely defined

Use only indices $s < i \leq r - s$:

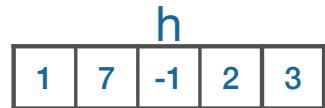
Size of next layer reduced by 2s



© 2020 YORAM SINGER 11

© 2020 YORAM SINGER 12

Padded Convolution



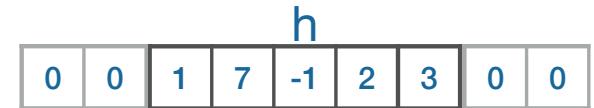
Define:

$$h_i = h_0 \text{ for } i \leq 0$$

$$h_i = h_{r+1} \text{ for } i > r$$

Output defined for $1 \leq i \leq r$

Padded Convolution



Define:

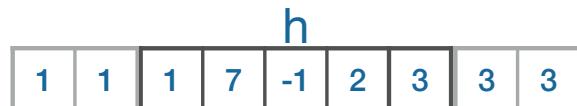
$$h_i = h_0 \text{ for } i \leq 0$$

$$\text{Zero padding: } h_{r+1} = h_0 = 0$$

$$h_i = h_{r+1} \text{ for } i > r$$

Output defined for $1 \leq i \leq r$

Padded Convolution



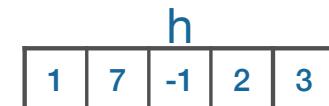
Define:

$$h_i = h_0 \text{ for } i \leq 0$$

$$h_i = h_{r+1} \text{ for } i > r \quad \text{Extension padding: } h_0 = h_1 \text{ and } h_{r+1} = h_r$$

Output defined for $1 \leq i \leq r$

Padded Convolution



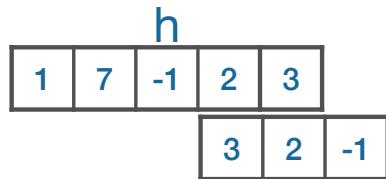
Define:

$$h_i = h_0 \text{ for } i \leq 0$$

$$h_i = h_{r+1} \text{ for } i > r$$

Output defined for $1 \leq i \leq r$

Padded Convolution

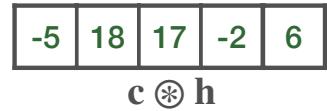


Define:

$$h_i = h_0 \text{ for } i \leq 0$$

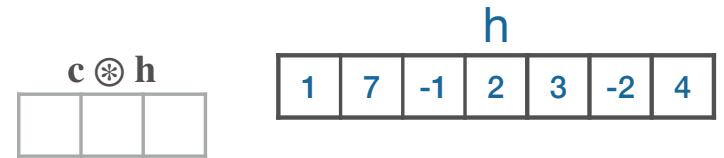
$$h_i = h_{r+1} \text{ for } i > r$$

Output defined for $1 \leq i \leq r$



© 2020 YORAM SINGER 13

Striding

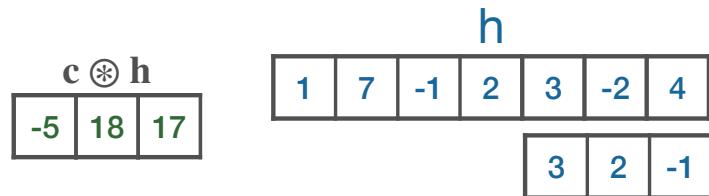


Reduce size of next layer through sub-sampling

$$h_i^{r+1} = \sigma\left(\sum_{j=-s}^s c_j h_{\mu(i)-j}^r\right) \text{ where } \mu(i) = b(i-1) + 1$$

© 2020 YORAM SINGER 14

Striding

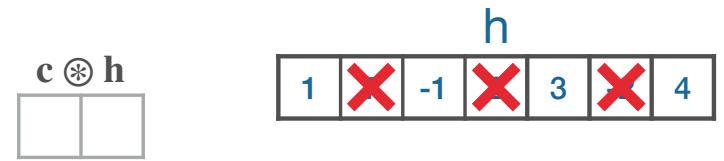


Reduce size of next layer through sub-sampling

$$h_i^{r+1} = \sigma\left(\sum_{j=-s}^s c_j h_{\mu(i)-j}^r\right) \text{ where } \mu(i) = b(i-1) + 1$$

© 2020 YORAM SINGER 14

Sub-Sampling (Dilation)

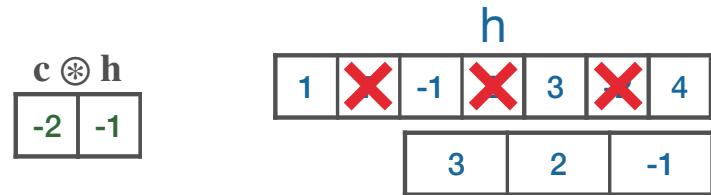


Reduce size of next layer through sub-sampling

$$h_i^{r+1} = \sigma\left(\sum_{j=-s}^s c_j h_{i-\mu(j)}^r\right) \text{ where } \mu(i) = b j$$

© 2020 YORAM SINGER 15

Sub-Sampling (Dilation)

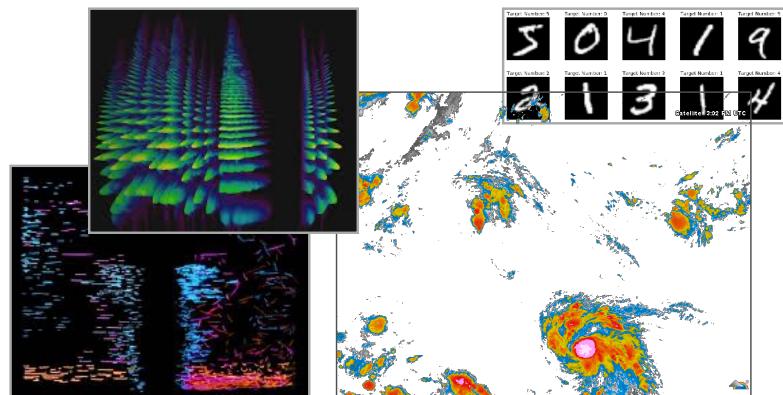


Reduce size of next layer through sub-sampling

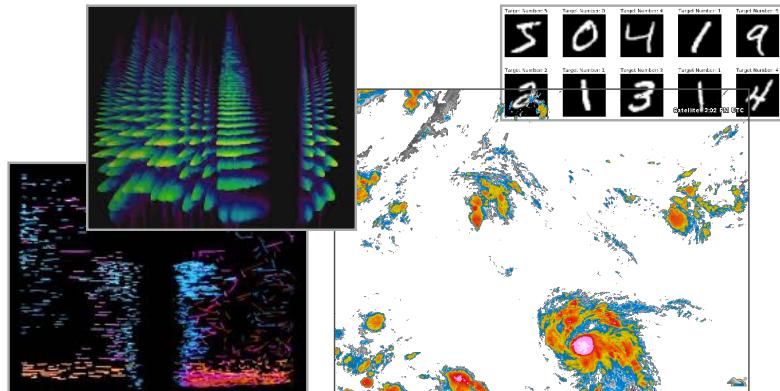
$$h_i^{r+1} = \sigma \left(\sum_{j=-s}^s c_j h_{i-\mu(j)}^r \right) \text{ where } \mu(i) = b_j$$

© 2020 YORAM SINGER 15

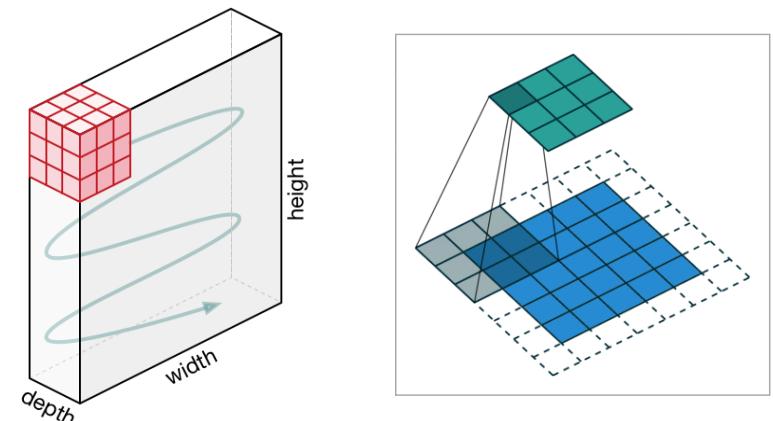
Multidimensional Convolutions



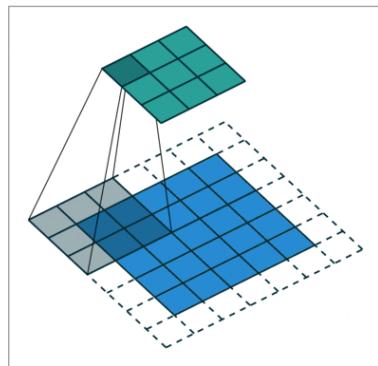
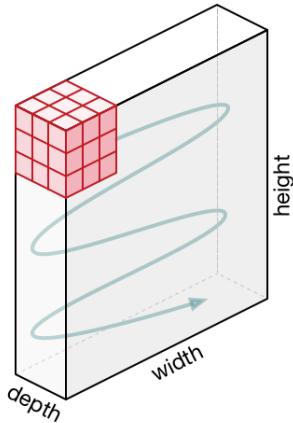
Multidimensional Convolutions



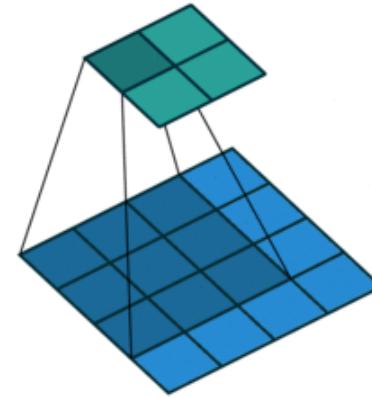
2D & 3D Convolutions



2D & 3D Convolutions

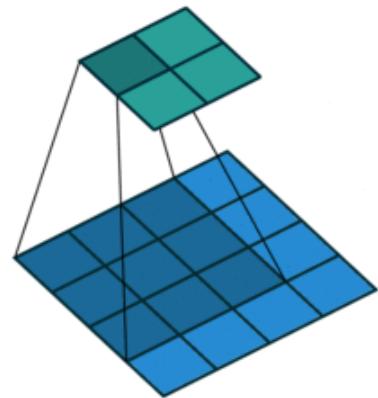


© 2020 YORAM SINGER 17



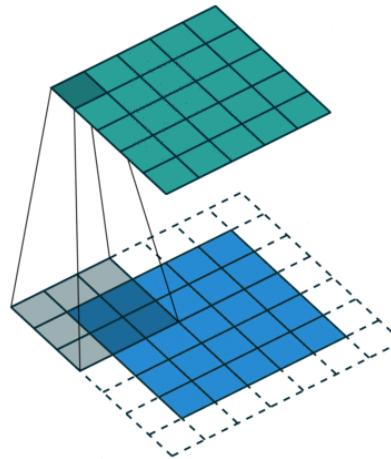
No Padding
No Striding
 $c \in \mathbb{R}^{3 \times 3}$

© 2020 YORAM SINGER 18



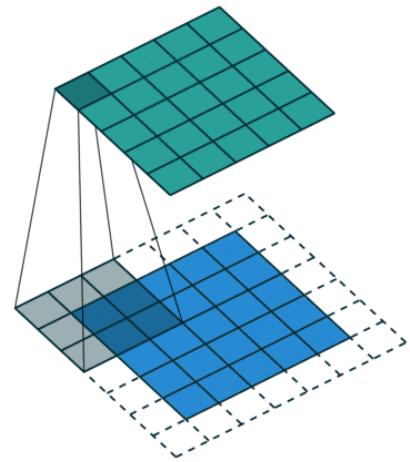
No Padding
No Striding
 $c \in \mathbb{R}^{3 \times 3}$

© 2020 YORAM SINGER 18



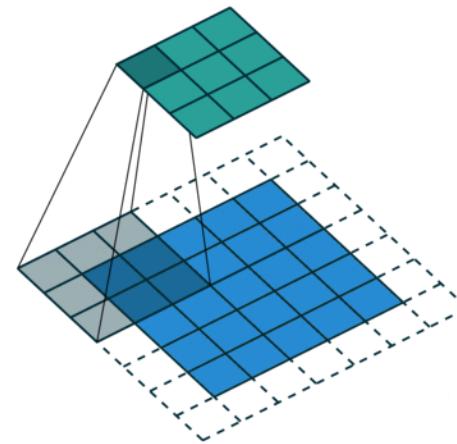
Padding
No Striding
 $c \in \mathbb{R}^{3 \times 3}$

© 2020 YORAM SINGER 19



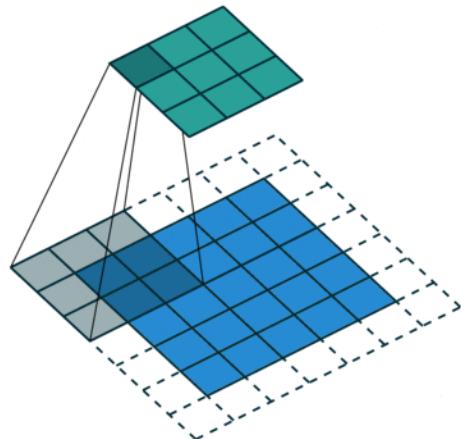
Padding
No Striding
 $c \in \mathbf{R}^{3 \times 3}$

© 2020 YORAM SINGER 19



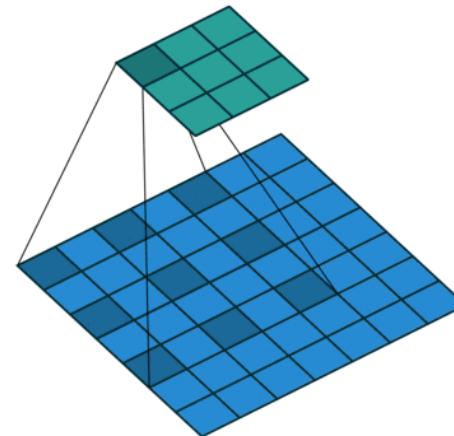
Padding
Striding
 $c \in \mathbf{R}^{3 \times 3}$

© 2020 YORAM SINGER 20



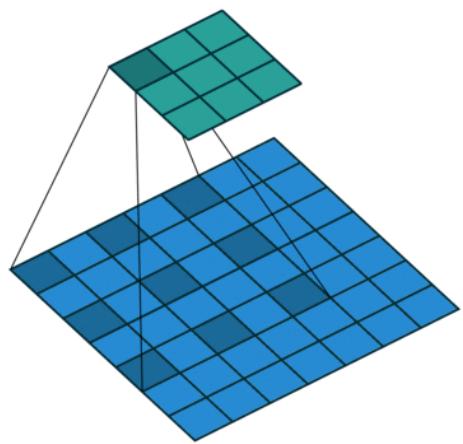
Padding
Striding
 $c \in \mathbf{R}^{3 \times 3}$

© 2020 YORAM SINGER 20



No Padding
Dilation
 $c \in \mathbf{R}^{3 \times 3}$

© 2020 YORAM SINGER 21

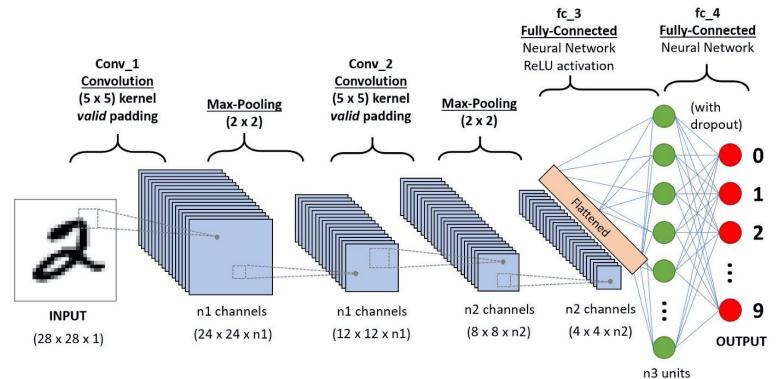


No Padding
Dilation
 $c \in \mathbb{R}^{3 \times 3}$

© 2020 YORAM SINGER 21

LeNet

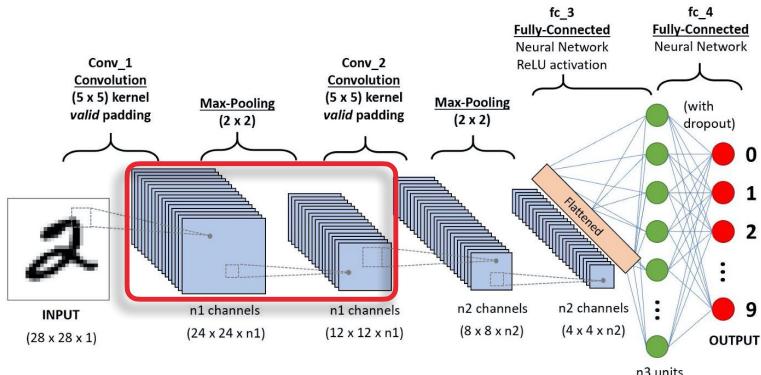
LeCun, Boser, Denker, Henderson, Howard, Hubbard, Jackel
Backpropagation applied to handwritten zip code recognition.
Neural Computation, 1989



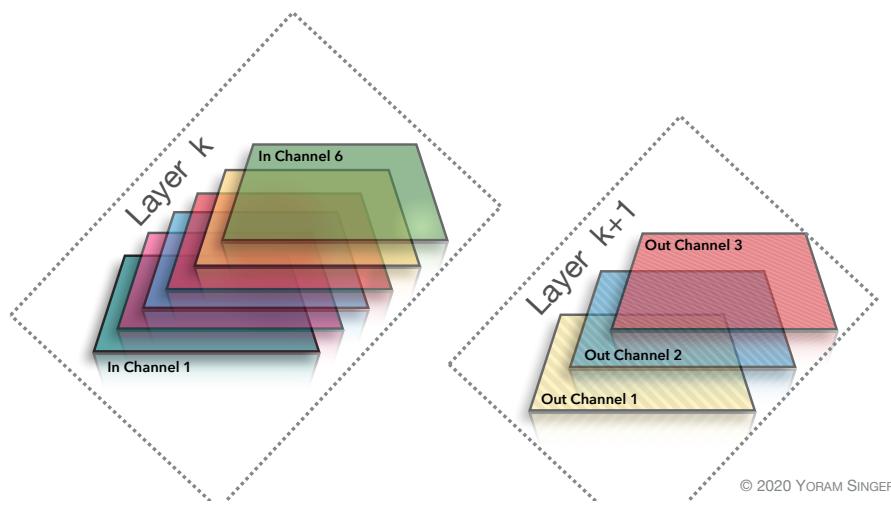
© 2020 YORAM SINGER 22

LeNet

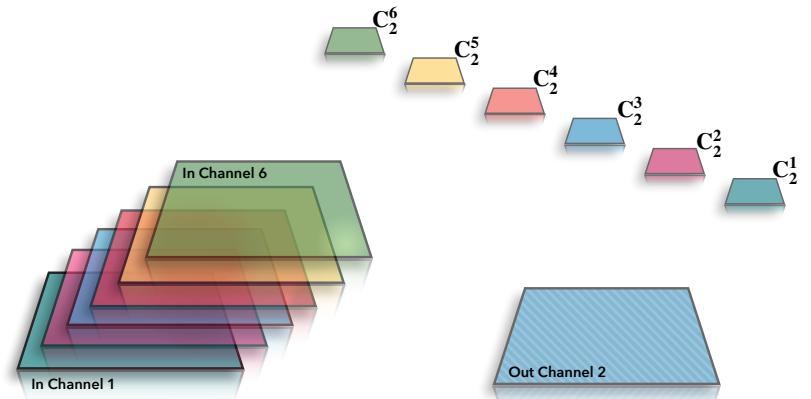
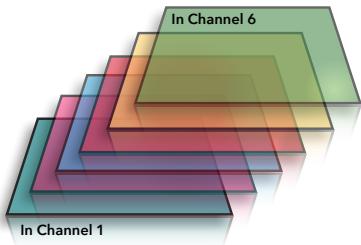
LeCun, Boser, Denker, Henderson, Howard, Hubbard, Jackel
Backpropagation applied to handwritten zip code recognition.
Neural Computation, 1989



© 2020 YORAM SINGER 22



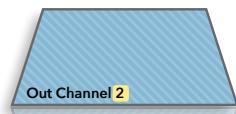
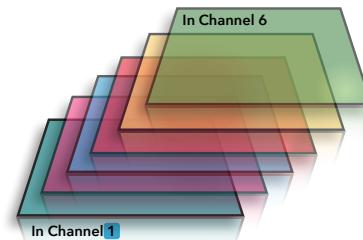
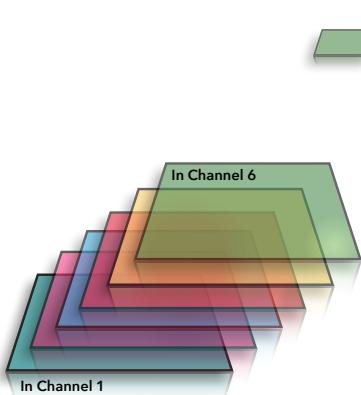
© 2020 YORAM SINGER 23



© 2020 YORAM SINGER 23

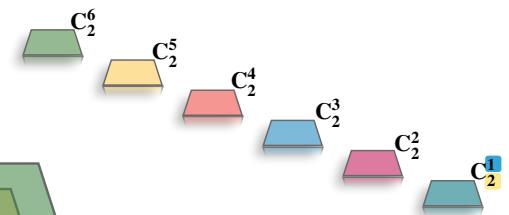
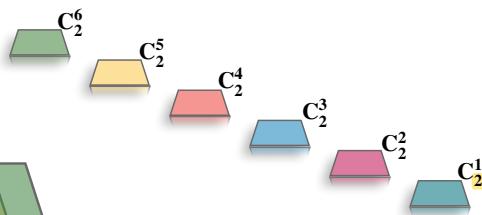


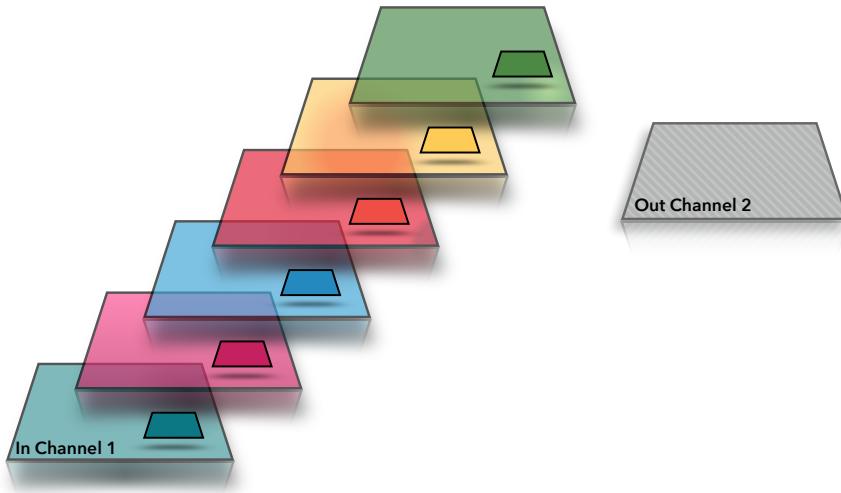
© 2020 YORAM SINGER 23



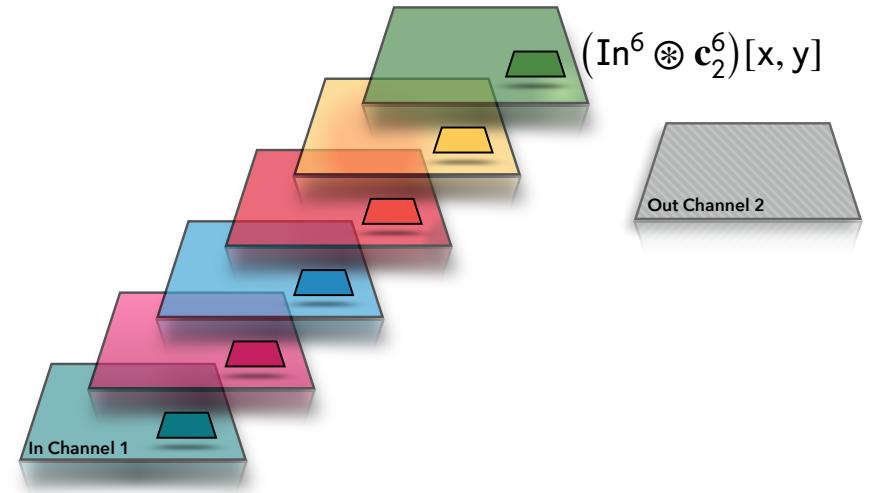
© 2020 YORAM SINGER 23

© 2020 YORAM SINGER 23

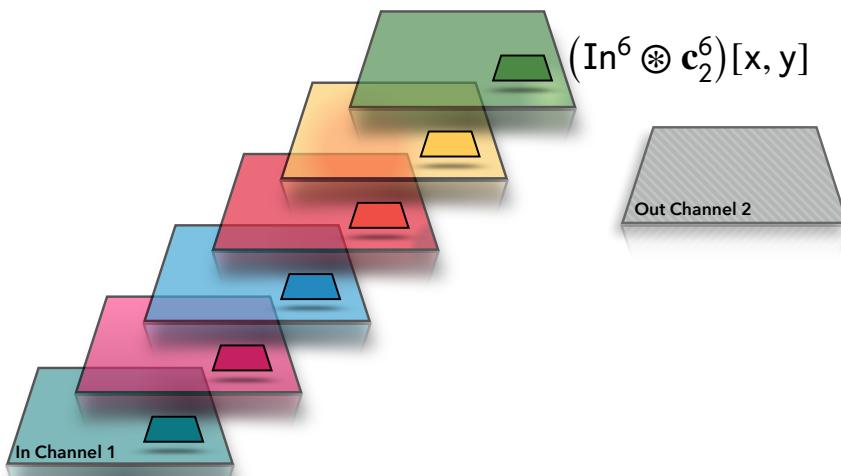




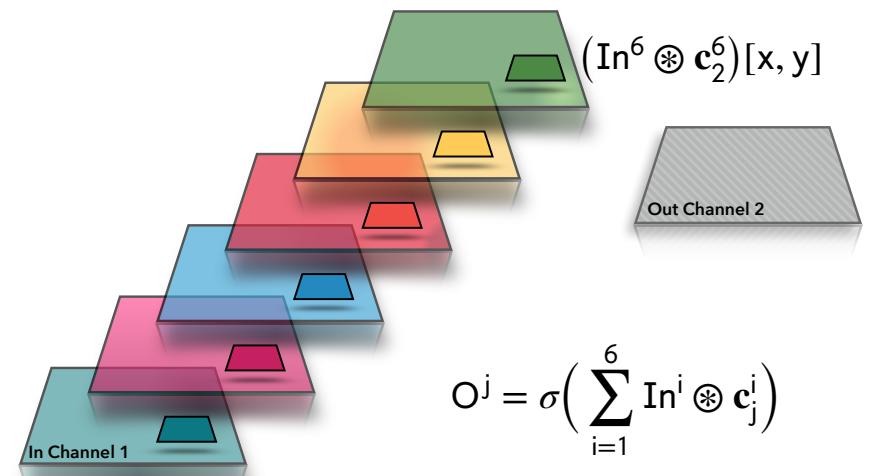
© 2020 YORAM SINGER 24



© 2020 YORAM SINGER 24

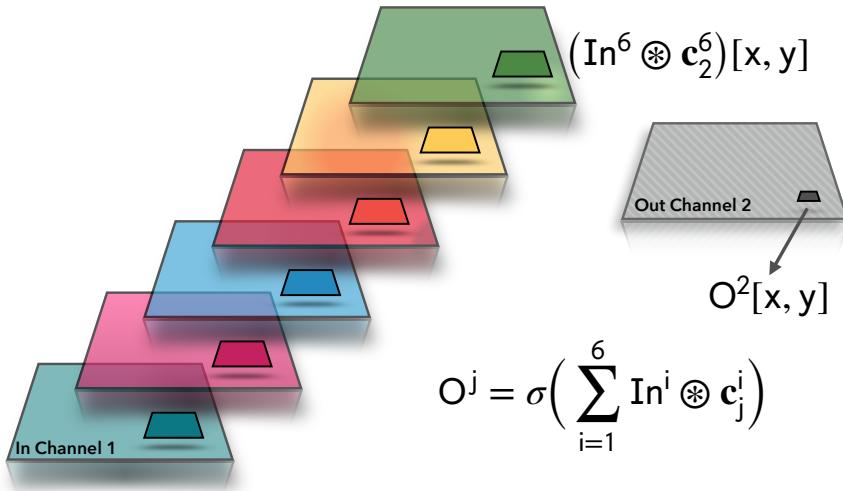


© 2020 YORAM SINGER 24



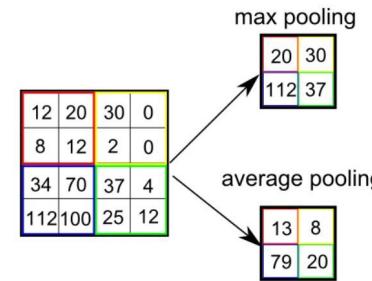
$$O^j = \sigma \left(\sum_{i=1}^6 I^{n^i} \otimes c_j^i \right)$$

© 2020 YORAM SINGER 24



© 2020 YORAM SINGER 24

2D Pooling & Averaging



2xx2 Max-Pooling reduces #neurons by 4

Tool for building a pyramid of layers

Channels flattened layer d-2 or d-3

Two-three fully connected layers

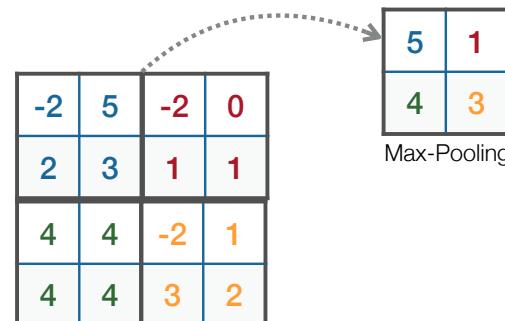
© 2020 YORAM SINGER 25

2D Max & Avg Pooling

-2	5	-2	0
2	3	1	1
4	4	-2	1
4	4	3	2

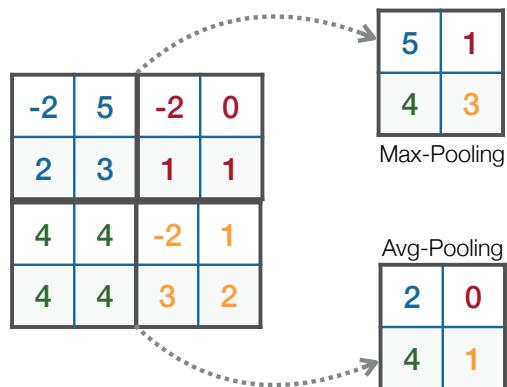
© 2020 YORAM SINGER 26

2D Max & Avg Pooling



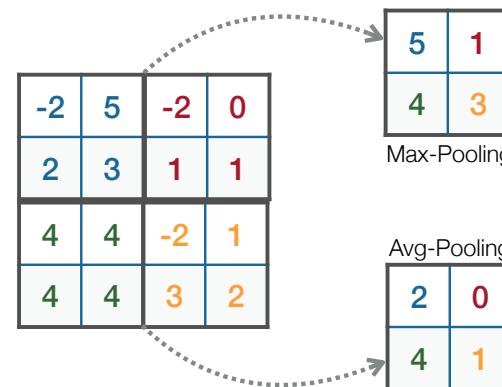
© 2020 YORAM SINGER 26

2D Max & Avg Pooling



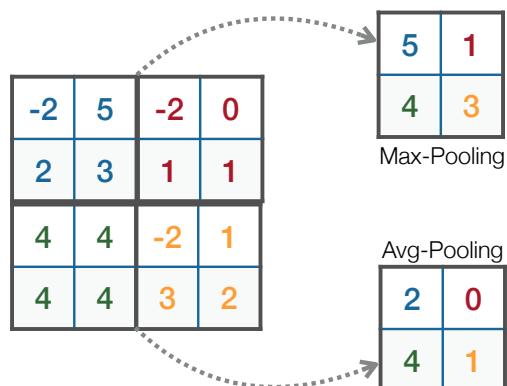
© 2020 YORAM SINGER 26

2D Max & Avg Pooling



© 2020 YORAM SINGER 26

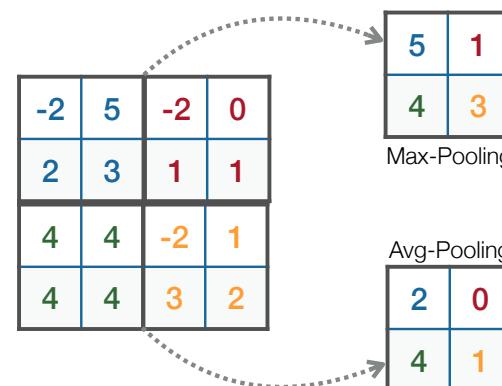
2D Max & Avg Pooling



- ▶ Pooling:
Reduces #neurons
Forms pyramid of layers
- ▶ Channels flattened to vector before top layers

© 2020 YORAM SINGER 26

2D Max & Avg Pooling



- ▶ Pooling:
Reduces #neurons
Forms pyramid of layers
- ▶ Channels flattened to vector before top layers
- ▶ Top layers fully connected

© 2020 YORAM SINGER 26

Code

1. Dataset: CIFAR10 for image classification, more difficult than MNIST for OCR
2. Notebooks: cifar10_mlp & cifar10_conv
3. Uses PyTorch with autograd: no need to calculate gradients by hand
4. Main loop written explicitly — examine carefully
5. Load and read notebook once before running cell by cell