



M2177.003100

Deep Learning

[7: Convolutional Neural Nets (Part 2)]

Electrical and Computer Engineering
Seoul National University

© 2019 Sungroh Yoon. this material is for educational uses only. some contents are based on the material provided by other paper/book authors and may be copyrighted by them.

(last compiled at 18:15:00 on 2019/10/13)

Outline

More on Convolution

- Backprop over Convolution

- 1×1 Convolution

- Transposed Convolution

- Other Types

Summary

References

- *Deep Learning* by Goodfellow, Bengio and Courville [▶ Link](#)
 - ▶ Chapter 9
- online resources:
 - ▶ *Deep Learning Specialization (coursera)* [▶ Link](#)
 - ▶ *Stanford CS231n: CNN for Visual Recognition* [▶ Link](#)
 - ▶ *Dive into Deep Learning (14.10 Transposed Convolution)* [▶ Link](#)
 - ▶ *Convolution Arithmetic* [▶ Link 1](#) [▶ Link 2](#) [▶ arXiv](#)
 - ▶ *Kunlun Bai's Blog on Convolution Types* [▶ Link](#)
- note:
 - ▶ you should [open this file in Adobe Acrobat](#) to see animated images (other types of pdf readers will not work)

Outline

More on Convolution

- Backprop over Convolution

- 1×1 Convolution

- Transposed Convolution

- Other Types

Summary

Recall: convolution

- with kernel flipping
 - ▶ commutative

$$\begin{aligned} Z(i, j) &= (K * V)(i, j) = \sum_m \sum_n \underbrace{K(m, n)}_{\text{kernel}} \underbrace{V(i - m, j - n)}_{\text{input volume}} \\ &= \sum_m \sum_n K(i - m, j - n) V(m, n) \\ &= (V * K)(i, j) \end{aligned}$$

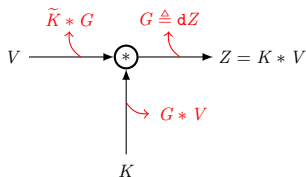
- without kernel flipping
 - ▶ not commutative

$$Z(i, j) = (K * V)(i, j) = \sum_m \sum_n \underbrace{K(m, n)}_{\text{kernel}} \underbrace{V(i + m, j + n)}_{\text{input volume}}$$

- ▶ we stick to this definition of convolution

Backprop over convolution

- forward



$$K * V = Z$$

$$C\mathbf{v} = \mathbf{z}$$

- backward

$$dV = \tilde{K} * G$$

$$d\mathbf{v} = C^T \mathbf{g}$$

$$dK = G * V$$

- ▶ $\mathbf{v}, \mathbf{z}, \mathbf{g}$: flattened versions of V, Z, G , respectively
- ▶ C : matrix representation of convolution
 - ▷ C^T : leads to _____ convolution
- ▶ \tilde{K} : flipped version of kernel K

Running example

- $K * V = Z$

- $(k, m, s, p)^1 = (2, 3, 1, 0)$

$$\begin{array}{|c|c|} \hline k_{11} & k_{12} \\ \hline k_{21} & k_{22} \\ \hline \end{array} * \begin{array}{|c|c|c|} \hline v_{11} & v_{12} & v_{13} \\ \hline v_{21} & v_{22} & v_{23} \\ \hline v_{31} & v_{32} & v_{33} \\ \hline \end{array} = \begin{array}{|c|c|} \hline z_{11} & z_{12} \\ \hline z_{21} & z_{22} \\ \hline \end{array}$$

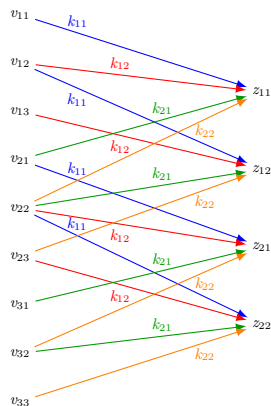
► result

$$z_{11} = k_{11}v_{11} + k_{12}v_{12} + k_{21}v_{21} + k_{22}v_{22}$$

$$z_{12} = k_{11}v_{12} + k_{12}v_{13} + k_{21}v_{22} + k_{22}v_{23}$$

$$z_{21} = k_{11}v_{21} + k_{12}v_{22} + k_{21}v_{31} + k_{22}v_{32}$$

$$z_{22} = k_{11}v_{22} + k_{12}v_{23} + k_{21}v_{32} + k_{22}v_{33}$$



¹sizes of kernel, input, stride, padding, respectively

Convolution as a matrix operation

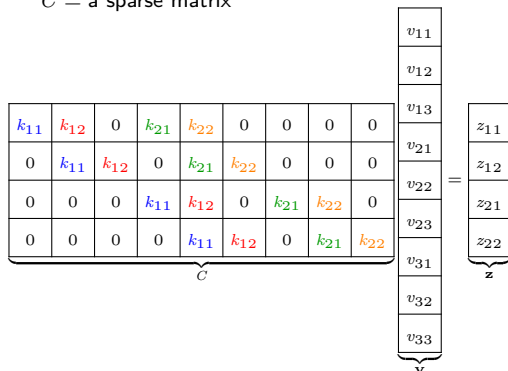
$$= \mathbf{z}$$

where

$$\mathbf{v} = (v_{11}, v_{12}, v_{13}, v_{21}, v_{22}, v_{23}, v_{31}, v_{32}, v_{33})$$

$$\mathbf{z} = (z_{11}, z_{12}, z_{21}, z_{22})$$

C = a sparse matrix

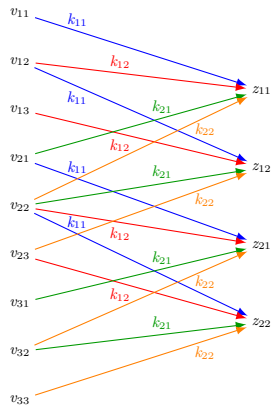


$$k_{11} v_{11} + k_{12} v_{12} + k_{21} v_{21} + k_{22} v_{22} = z_{11}$$

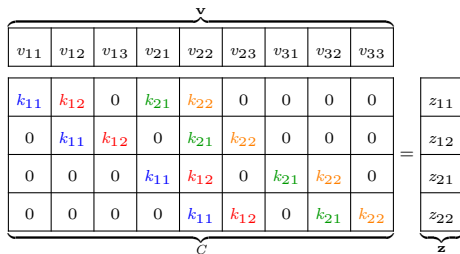
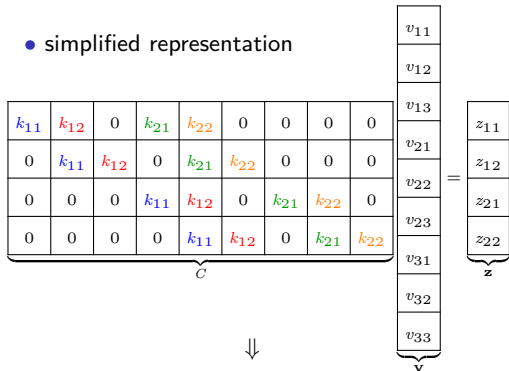
$$k_{11} v_{12} + k_{12} v_{13} + k_{21} v_{22} + k_{22} v_{23} = z_{12}$$

$$k_{11} v_{21} + k_{12} v_{22} + k_{21} v_{31} + k_{22} v_{32} = z_{21}$$

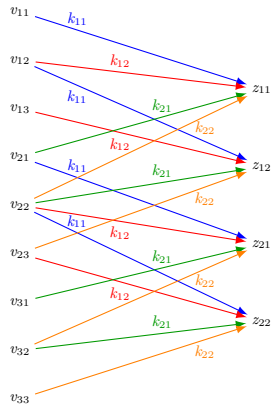
$$k_{11} v_{22} + k_{12} v_{23} + k_{21} v_{32} + k_{22} v_{33} = z_{22}$$



• simplified representation



$$\begin{aligned}
 k_{11} v_{11} + k_{12} v_{12} + k_{21} v_{21} + k_{22} v_{22} &= z_{11} \\
 k_{11} v_{12} + k_{12} v_{13} + k_{21} v_{22} + k_{22} v_{23} &= z_{12} \\
 k_{11} v_{21} + k_{12} v_{22} + k_{21} v_{31} + k_{22} v_{32} &= z_{21} \\
 k_{11} v_{22} + k_{12} v_{23} + k_{21} v_{32} + k_{22} v_{33} &= z_{22}
 \end{aligned}$$



$$dV = \tilde{K} * G$$

$$dv_{11} = k_{11}g_{11}$$

$$dv_{12} = k_{12}g_{11} + k_{11}g_{12}$$

$$dv_{13} = k_{12}g_{12}$$

$$dv_{21} = k_{21}g_{11} + k_{11}g_{21}$$

$$dv_{22} = k_{22}g_{11} + k_{21}g_{12} + k_{12}g_{21} + k_{11}g_{22}$$

$$dv_{23} = k_{22}g_{12} + k_{12}g_{22}$$

$$dv_{31} = k_{21}g_{21}$$

$$dv_{32} = k_{22}g_{21} + k_{21}g_{22}$$

$$dv_{33} = k_{22}g_{22}$$

dv_{11}	dv_{12}	dv_{13}
dv_{21}	dv_{22}	dv_{23}
dv_{31}	dv_{32}	dv_{33}

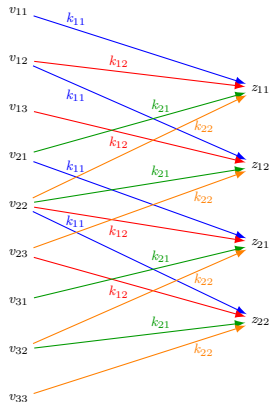
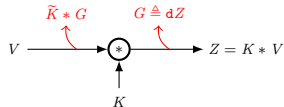
$$=$$

k_{22}	k_{21}
k_{12}	k_{11}

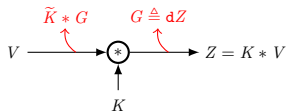
$$*$$

0	0	0	0
0	g_{11}	g_{12}	0
0	g_{21}	g_{22}	0
0	0	0	0

\tilde{K} : flipped K (padded) $G=dZ$



- matrix representation (reveals “_____ convolution”)



backward

$$dV = \tilde{K} * G$$

$$d\mathbf{v} = C^T \mathbf{g}$$

forward

$$K * V = Z$$

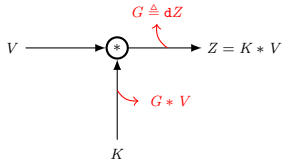
$$C\mathbf{v} = \mathbf{z}$$

\mathbf{v}								
v_{11}	v_{12}	v_{13}	v_{21}	v_{22}	v_{23}	v_{31}	v_{32}	v_{33}
k_{11}	k_{12}	0	k_{21}	k_{22}	0	0	0	0
0	k_{11}	k_{12}	0	k_{21}	k_{22}	0	0	0
0	0	0	k_{11}	k_{12}	0	k_{21}	k_{22}	0
0	0	0	0	k_{11}	k_{12}	0	k_{21}	k_{22}
C								

\mathbf{z}
z_{11}
z_{12}
z_{21}
z_{22}

\mathbf{g}			
g_{11}	g_{12}	g_{21}	g_{22}
dv_{11}	k_{11}	0	0
dv_{12}	k_{12}	k_{11}	0
dv_{13}	0	k_{12}	0
dv_{21}	k_{21}	0	k_{11}
dv_{22}	k_{22}	k_{21}	k_{12}
dv_{23}	0	k_{22}	0
dv_{31}	0	0	k_{21}
dv_{32}	0	0	k_{22}
dv_{33}	0	0	k_{22}
C^T			

$$dK = G * V$$



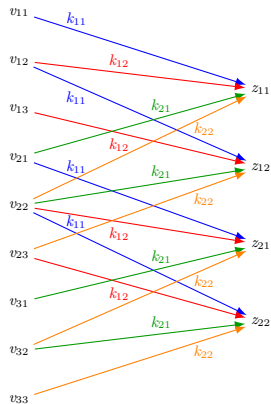
$$dk_{11} = g_{11}v_{11} + g_{12}v_{12} + g_{21}v_{21} + g_{22}v_{22}$$

$$dk_{12} = g_{11}v_{12} + g_{12}v_{13} + g_{21}v_{22} + g_{22}v_{23}$$

$$dk_{21} = g_{11}v_{21} + g_{12}v_{22} + g_{21}v_{31} + g_{22}v_{32}$$

$$dk_{22} = g_{11}v_{22} + g_{12}v_{23} + g_{21}v_{32} + g_{22}v_{33}$$

$$\underbrace{\begin{bmatrix} dk_{11} & dk_{12} \\ dk_{21} & dk_{22} \end{bmatrix}}_{dK} = \underbrace{\begin{bmatrix} g_{11} & g_{12} \\ g_{21} & g_{22} \end{bmatrix}}_G * \underbrace{\begin{bmatrix} v_{11} & v_{12} & v_{13} \\ v_{21} & v_{22} & v_{23} \\ v_{31} & v_{32} & v_{33} \end{bmatrix}}_V$$



Outline

More on Convolution

Backprop over Convolution

1×1 Convolution

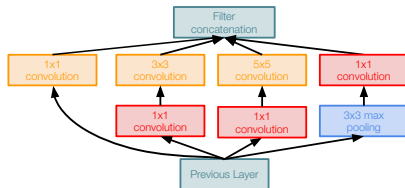
Transposed Convolution

Other Types

Summary

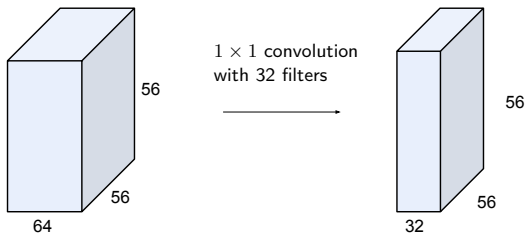
1×1 convolution

- aka **pointwise** convolution
- widely used for _____ adjustment
e.g. inception module in GoogLeNet



(source: Szegedy et al., 2014)

- example:



(source: cs231n)

- ▶ each filter: size $1 \times 1 \times 64$ (performs a 64-dim dot product)

1×1 convolution on volumes

6	2	1	1
9	4	5	6
2	7	5	3
8	2	3	5

4×4

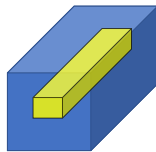
*

2

=

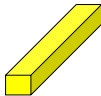
12	4	2	2
18	8	10	12
4	14	10	6
16	4	6	10

4×4



$4 \times 4 \times 64$

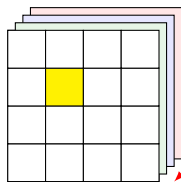
*



$1 \times 1 \times 64$

C filters

=



4×4

$4 \times 4 \times C$

depth C

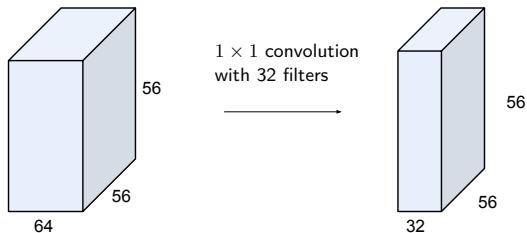
(source: coursera)

- nonlinearity (e.g. ReLU) can follow 1×1 convolution \rightarrow “network in network”

Depth adjustment

- 1×1 convolution: widely used for depth adjustment

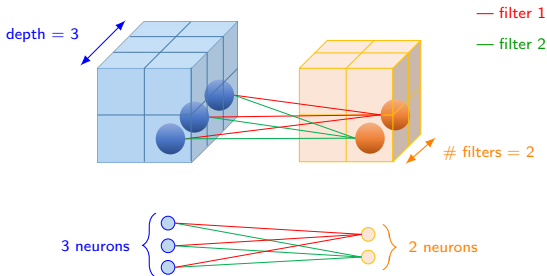
e.g.



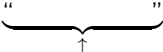
(source: cs231n)

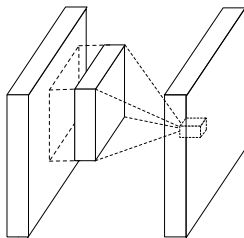
- ▶ each filter: size $1 \times 1 \times 64$ (performs a 64-dim dot product)
 - ▶ preserves spatial dimensions and reduces depth
 - ▶ projects depth to ____ dimension (combination of feature maps)
- in general
 - ▶ we can reduce/maintain/increase depth using 1×1 convolution

- a set of 1×1 conv filters: can be interpreted as forming an _____
 - ▶ **input dimension** of this FC layer
= _____ of the input volume to 1×1 conv filters
 - ▶ **output dimension** of this FC layer
= _____ of 1×1 conv filters
- example: $2 \times 2 \times 3$ volume applied to two $1 \times 1 \times 3$ filters
 - ▶ $2 \times 2 = 4$ applications of the FC layer that maps 3 neurons to 2 neurons

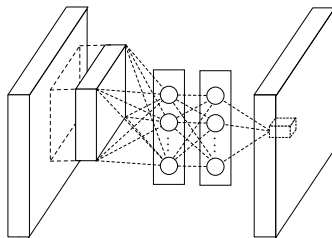


Network in network

- Mlpconv layer with “” within each conv layer
 - composed of FC layer (with 1×1 conv) + nonlinearity
- ▶ can compute more **abstract features** for local patches
- ▶ precursor to GoogLeNet and ResNet “bottleneck” layers



(a) linear convolution layer



(b) Mlpconv layer

(source: Lin et al., 2014)

Outline

More on Convolution

Backprop over Convolution

1×1 Convolution

Transposed Convolution

Other Types

Summary

Motivation: upsampling

- desire to use a transform going in opposite direction of normal convolution

e.g. _____ layer of a convolutional autoencoder
project feature maps to a higher-dim space

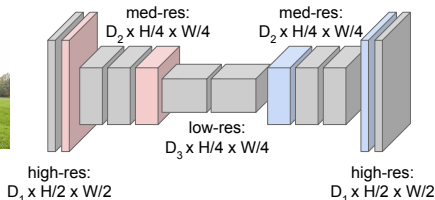
- example in cv: semantic segmentation

downsampling:
pooling, strided
convolution



input:
 $3 \times H \times W$

design network as a bunch of convolutional layers, with
downsampling and **upsampling** inside the network!



upsampling:
???



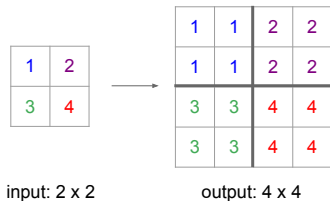
predictions:
 $H \times W$

(source: cs231n)

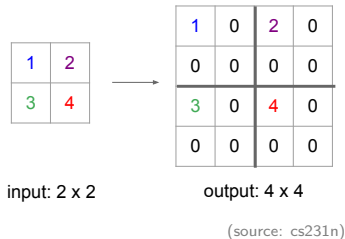
- ▶ downsampling: convolution + pooling
- ▶ how to do upsampling?

Rule-based upsampling

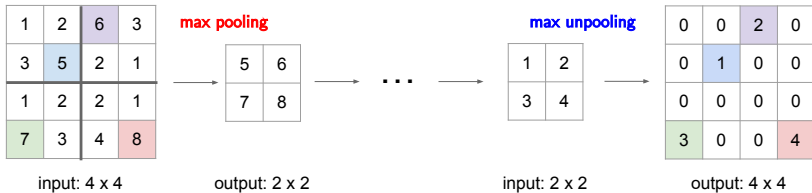
- nearest neighbor



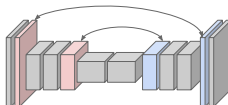
- “bed of nails”



- max unpooling: remember ____ positions from pooling layer



(source: cs231n)



corresponding pairs of
downsampling and
upsampling layers

Transposed convolution

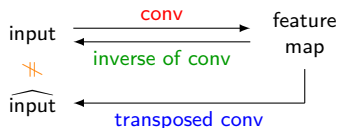
- _____ upsampling: smaller map \rightarrow larger map
 - c.f. **regular convolution**: downsampling (larger map \rightarrow smaller map)
 - ▶ does not perform the **deconvolution** (*i.e.* **inverse of convolution**)²
 - ▶ *aka* fractionally strided³ conv, upconv, backward strided conv

- example

(source: Theano tutorial)

- ▶ 3×3 filter
 - ▶ 2×2 input $\rightarrow 5 \times 5$ output

- note:

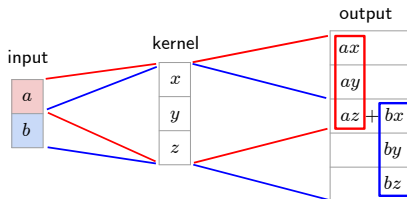


²thus, transposed convolution is sometimes (inappropriately) called *deconvolution* but is different from the deconvolution in engineering and mathematics

³stride: gives ratio between movement in output and input (in pixels/out pixels)

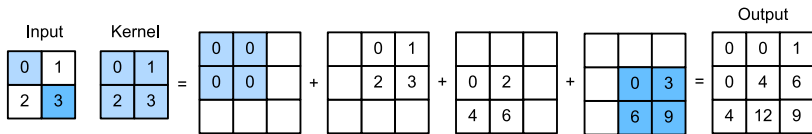
more examples:

- 1D transposed convolution (3×1 kernel)
 - ▶ output contains copies of kernel weighted by input
 - ▶ overlaps are _____



(source: cs231n)

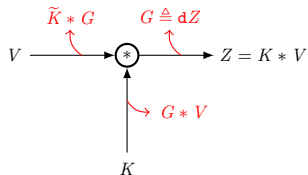
- 2D transposed convolution (2×2 kernel)



(source: dive into DL)

More formally

- recall: backprop over convolution



- ▶ forward

$$K * V = Z$$

$$C\mathbf{v} = \mathbf{z}$$

- ▶ backward

$$dV = \tilde{K} * G$$

$$d\mathbf{v} = C^\top \mathbf{g}$$

- transposed convolution

- ▶ defined as _____ pass of regular convolution with the same kernel

$$\text{TransposedConv}(K, X) = \tilde{K} * X$$

- ▶ more details:

▶ [Link 1](#)

▶ [Link 2](#)

- the kernel K defines a **convolution**
 - ▶ whose forward and backward passes:
 - ▷ computed by multiplying C and C^\top , respectively
- K also defines a **transposed convolution**
 - ▶ whose forward and backward passes:
 - ▷ computed by multiplying C^\top and $(C^\top)^\top = C$, respectively
- implementation of transposed convolution
 - ▶ implement with a **normal convolution** + (complicated) **zero-padding**
 - ▷ easy to understand but much less efficient
 - ▶ implement as _____ **of some convolution** wrt its input
 - ▷ usually how it is implemented in practice

Outline

More on Convolution

Backprop over Convolution

1×1 Convolution

Transposed Convolution

Other Types

Summary

Dilated convolution (atrous convolution)

- introduces another convolution parameter: dilation rate
 - ▶ defines a spacing between values in a filter
- e.g. 3×3 filter with dilation rate 2
 - ⇒ the same field of view as 5×5 filter (but only uses 9 parameters)
- effect: a _____ field of view at the same computational cost
 - ▶ real-time segmentation, speech synthesis
- example
 - ▶ 3×3 filter
 - ▶ dilation rate of 2
 - ▶ no zero-padding

(source: Theano tutorial)

Separable convolution

- two types
 - ▶ **spatially** separable conv
 - ▶ separable conv: popular (*e.g.* MobileNet, Xception)
 - **spatially** separable convolution
 - ▶ operates on 2D spatial dimensions (height and width)
 - ▶ decomposes a convolution into two separate operations
- i.e.* a 2D kernel \rightarrow two 1D kernels

$$\begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} \begin{bmatrix} -1 & 0 & 1 \end{bmatrix}$$

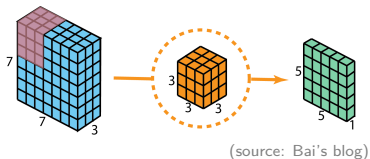
- ▶ applies each 1D kernel in turn [▶ example](#)
- \Rightarrow reduction in computation compared with regular convolution
- ▶ rarely used in deep learning (not every 2D kernel is decomposable)

- depthwise separable convolution: two-step process

1. depthwise convolution
2. pointwise (1×1) convolution

- e.g. $7 \times 7 \times 3$ image $\rightarrow 5 \times 5 \times 1$ feature map $[(k, m, s, p) = (3, 7, 1, 0)]$

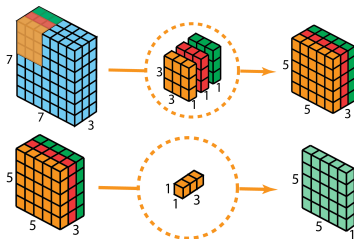
regular convolution



$$(3 \times 3 \times 3) \times (5 \times 5) = 675 \text{ mults}$$

___ computational gain for 1 filter!

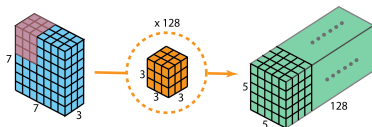
depthwise separable convolution



$$\begin{aligned}
 3 \times (3 \times 3 \times 1) \times (5 \times 5) &= 675 \\
 + \\
 (1 \times 1 \times 3) \times (5 \times 5) &= 75 \\
 \text{total: } 675 + 75 &= 750 \text{ mults}
 \end{aligned}$$

- how about multiple (e.g. 128) filters?

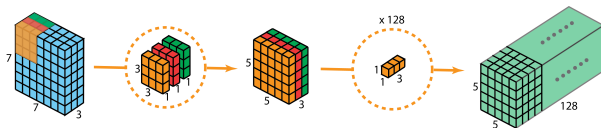
► **regular** convolution



(source: Bai's blog)

$$128 \times (3 \times 3 \times 3) \times (5 \times 5) = 128 \times 675 = 86,400 \text{ mults}$$

► **depthwise separable** convolution



$$3 \times (3 \times 3 \times 1) \times (5 \times 5) + 128 \times (1 \times 1 \times 3) \times (5 \times 5) = 9,600 \text{ mults}$$

⇒ huge computational gain! (only 11% of regular conv)

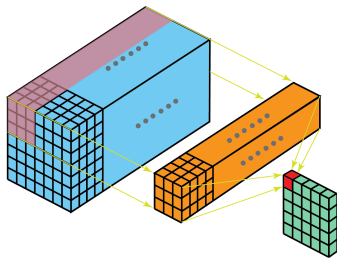
- significantly fewer kernel parameters ($128 \times 27 = 3,456$ vs $27 + 128 \times 3 = 411$)

⇒ reduced model _____ (problematic if not properly trained)

3D convolution

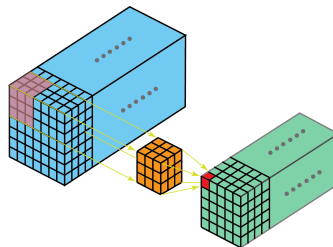
- regular (2D) convolution

- ▶ kernel depth = input depth
- ▶ kernel moves only in 2D (width, height)



- 3D convolution

- ▶ kernel depth < input depth
- ▶ kernel moves in all ____ (width, height, depth)

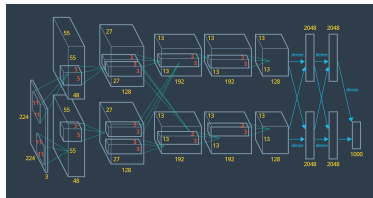


(source: Bai's blog)

- ▶ d -dim conv: describes spatial relationships of object in d -dim space

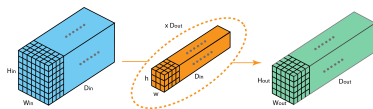
Grouped convolution

- first introduced in AlexNet (2012)
 - ▶ mainly due to limited _____
 - ▶ other advantages also exist
- ▶ details

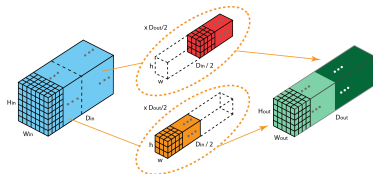


(source: yuchao.us)

- regular convolution



- grouped convolution (2 groups)



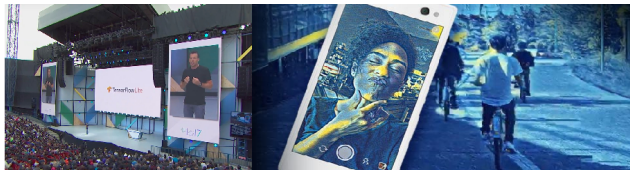
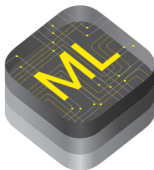
(source: Bai's blog)

Convolution by frequency domain conversion

- convolution: equivalent to the following
 1. convert both input/kernel to _____ domain using Fourier transform
 2. perform point-wise multiplication of the two signals
 3. convert back to time domain using an inverse Fourier transform
- for some problem sizes
 - ▶ this can be faster than the naïve implementation of discrete convolution

Remarks

- active areas of research
 - ▶ devising faster ways of performing convolution
 - ▶ approximate convolution without harming accuracy of the model
 - ▶ fast evaluation of forward propagation
 - in commercial setting
 - ▶ typically devote more resources to deployment of a net than its training
- ⇒ techniques that improve efficiency of only _____ prop are useful
- e.g.* TensorRT, TensorFlowLite, Core ML, Caffe2Go



Outline

More on Convolution

- Backprop over Convolution

- 1×1 Convolution

- Transposed Convolution

- Other Types

Summary

Summary

- backprop over convolution

- ▶ forward

$$K * V = Z$$

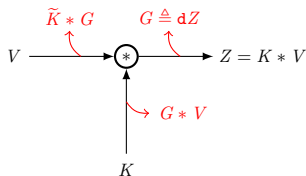
$$C\mathbf{v} = \mathbf{z}$$

- ▶ backward

$$dV = \tilde{K} * G$$

$$d\mathbf{v} = C^\top \mathbf{g}$$

$$dK = G * V$$



- various types of convolution operations exist (their fast inference: crucial)

- ▶ pointwise (1×1) convolution: for depth adjustment
 - ▶ transposed (fractionally strided) convolution: for learnable upsampling
 - ▶ dilated convolution, separable convolution: for efficiency (+ alpha)
 - ▶ 3D convolution, grouped convolution, and many others