

THE UNIVERSITY OF BRITISH COLUMBIA

Topics in AI (CPSC 532S): **Multimodal Learning with Vision, Language and Sound**

Lecture 4: Convolutional Neural Networks (Part 1)



Course Logistics

- Azure credits have been distributed (to those who asked for them)
- Assignment 1 was due yesterday
- Assignment 2 will be out today (on CNNs) and is due Wednsday next week (note, it will take computation time)

 Start thinking about **project** and forming groups (groups will be due in 1-2 weeks, proposal is in ~month)

- Instructions for using Azure at TA Office Hour 1-2pm today and on Piazza

Research **Projects**

Business

entrepreneur will take the best idea and run it into a ground."

Research

best project and run it into a ground." — Me

"A good entrepreneur can take a mediocre idea and make it great, a bad



Robert Herjavec (Canadian businessman and investor) Shark Tank

"A good idea can make a mediocre project great, a bad idea will take the

I want to **solve vision** / language / etc.



Marvin Minsky

- I want to **solve vision** / language / etc.
- I want to do X (e.g., image captioning)

 - Requires **forward thinking**, knowledge of the field
 - A sure way to get tenure
 - Difficult to do as the field matures

- This is excellent if X is something no one has done or thought about and is important (guaranteed success)



Rosland Picard, MIT — Affective Computing



- I want to **solve vision** / language / etc.
- I want to do X (e.g., image captioning)

think the **right way to solve** (or improve) X is Y

- More incremental; a lot of science is incremental ("standing on the shoulders of giants")
- **Retrospective:** compare existing approaches see why they work what is missing (guaranteed success)
- **Perspective:** come up with an idea or the insight that you truly believe and test it
- Requires through knowledge of the sub-field (lots of reading)
- Requires strong intuition and high level (intuitive) thinking
- Requires understanding of the mathematical tools and formulations to know what maybe possible
- Helps to bringing knowledge from other fields (field cross pollination)



Geoffrey Hinton







- I want to **solve vision** / language / etc.
- I want to do X (e.g., image captioning)
- I think the **right way to solve** (or improve) X is Y
- Mathematical formulation
- Implementation / engineering
- **Experimental** testing

On to todays lecture ...

Fully Connected Layer



Example: 200 x 200 image (small) x 40K hidden units

= ~ 2 Billion parameters (for one layer!)

Spatial correlations are generally local

Waste of resources + we don't have enough data to train networks this large





Locally Connected Layer



Example: 200 x 200 image (small)

Filter size: 10 x 10

= ~ 4 Million parameters

Locally Connected Layer



Example: 200 x 200 image (small)

Filter size: 10 x 10

= ~ 4 Million parameters

Stationarity — statistics is similar at different locations



Example: 200 x 200 image (small)

Filter size: 10×10

= ~ 4 Million parameters

= 100 parameters

Share the same parameters across the locations (assuming input is stationary)

* slide adopted from Marc'Aurelio Renzato











































$\begin{bmatrix} 0.11 & 0.11 & 0.11 \end{bmatrix}$ $\begin{bmatrix} 0.11 & 0.11 & 0.11 \\ 0.11 & 0.11 & 0.11 \end{bmatrix}$





Example: 200 x 200 image (small) x 40K hidden units

Filter size: 10 x 10

of filters: 20

= 2000 parameters

Learn multiple filters

32 x 32 x 3 image (note the image preserves spatial structure)



3 depth

32 x 32 x 3 **image**





$5 \times 5 \times 3$ filter

Convolve the filter with the image (i.e., "slide over the image spatially, computing dot products")



Convolutional Layer 32 x 32 x 3 image





Filters always extend the full depth of the input volume

5 x 5 x 3 filter

Convolve the filter with the image (i.e., "slide over the image spatially, computing dot products"





32 x 32 x 3 **image**





1 number: the result of taking a dot product between the filter and a small 5 x 5 x 3 part of the image

$$\mathbf{W}^T \mathbf{x} + b$$
, where $\mathbf{W}, \mathbf{x} \in \mathbb{R}^{75}$

How many **parameters** does the layer have? **76**

32 x 32 x 3 **image**





activation map
Convolutional Layer

32 x 32 x 3 **image**





activation map

Convolutional Layer





this results in the "new image" of size 28 x 28 x 6!



Convolutional Neural Network (ConvNet)







What filters do networks learn?



Layer 1





[Zeiler and Fergus, 2013]



What filters do networks learn?



[Zeiler and Fergus, 2013]



32 x 32 x 3 **image**





activation map



7 width



7 x 7 input image (spatially) 3 x 3 filter

7 height



7 width



7 x 7 input image (spatially) 3 x 3 filter

7 height



7 width



7 x 7 input image (spatially) 3 x 3 filter

7 height



7 width



7 x 7 input image (spatially) 3 x 3 filter

7 height



7 width



7 x 7 input image (spatially) 3 x 3 filter

=> **5** x **5** output

7 height



7 width



7 x 7 input image (spatially) 3 x 3 filter (applied with stride 2)

7 height



7 width



7 x 7 input image (spatially) 3 x 3 filter (applied with stride 2)

7 height



7 width



7 x 7 input image (spatially) 3 x 3 filter (applied with stride 2)

=> **3 x 3 output**

7 height



7 width



7 x 7 input image (spatially) 3 x 3 filter (applied with stride 3)

7 height

Does not fit! **Cannot apply** 3 x 3 filter on 7 x 7 image with stride 3



N width



N x N input image (spatially) F x F filter

Output size: (N-F) / stride + 1

N height

Example: N = 7, F = 3

stride $1 = \frac{(7-3)}{1+1} = 5$ stride 2 = (7-3)/2 + 1 = 3stride 3 => (7-3)/3+1 = **2.33**



Convolutional Layer: Border padding

7 width



7 x 7 input image (spatially)3 x 3 filter(applied with stride 1)

pad with 1 pixel border

7 height

Convolutional Layer: Border padding

7 width



7 x 7 input image (spatially)3 x 3 filter(applied with stride 1)

pad with 1 pixel border

Output size: 7 x 7

7 height

Convolutional Layer: Border padding

7 width



7 x 7 input image (spatially)3 x 3 filter(applied with **stride 3**)

pad with 1 pixel border

7 height

Example: N = 7, F = 3

stride 1 => (9-3)/1+1 = 7stride 2 => (9-3)/2+1 = 4stride 3 => (9-3)/3+1 = 3

Convolutional Neural Network (ConvNet)

With padding we can achieve no shrinking (32 -> 28 -> 24); shrinking quickly (which happens with larger filters) doesn't work well in practice





Convolutional Layer: 1x1 convolutions

56 x 56 x 64 **image**



Accepts a volume of size: $W_i \times H_i \times D_i$

- Accepts a volume of size: $W_i \times H_i \times D_i$ Requires hyperparameters:

 - Number of filters: K (for typical networks $K \in \{32, 64, 128, 256, 512\}$) - Spatial extent of filters: F (for a typical networks $F \in \{1, 3, 5, ...\}$) - Stride of application: S (for a typical network $S \in \{1, 2\}$) - Zero padding: P (for a typical network $P \in \{0, 1, 2\}$)

- Accepts a volume of size: $W_i \times H_i \times D_i$ Requires hyperparameters:

 - Number of filters: K (for typical networks $K \in \{32, 64, 128, 256, 512\}$) - Spatial extent of filters: F (for a typical networks $F \in \{1, 3, 5, ...\}$) - Stride of application: S (for a typical network $S \in \{1, 2\}$) - Zero padding: P (for a typical network $P \in \{0, 1, 2\}$)
- Produces a volume of size: $W_o \times H_o \times D_o$

- Accepts a volume of size: $W_i \times H_i \times D_i$ Requires hyperparameters:

 - Number of filters: K (for typical networks $K \in \{32, 64, 128, 256, 512\}$) - Spatial extent of filters: F (for a typical networks $F \in \{1, 3, 5, ...\}$) - Stride of application: S (for a typical network $S \in \{1, 2\}$) - Zero padding: P (for a typical network $P \in \{0, 1, 2\}$)
- Produces a volume of size: $W_o \times H_o \times D_o$
 - $W_o = (W_i F + 2P)/S + 1$ $H_o = (H_i F + 2P)/S + 1$

 $D_o = K$

- Accepts a volume of size: $W_i \times H_i \times D_i$ Requires hyperparameters:

 - Number of filters: K (for typical networks $K \in \{32, 64, 128, 256, 512\}$) - Spatial extent of filters: F (for a typical networks $F \in \{1, 3, 5, ...\}$) - Stride of application: S (for a typical network $S \in \{1, 2\}$) - Zero padding: P (for a typical network $P \in \{0, 1, 2\}$)
- Produces a volume of size: $W_o \times H_o \times D_o$

$$W_o = (W_i - F + 2P)/S + 1$$

Number of total learnable parameters: $(F \times F \times D_i) \times K + K$

- $H_o = (H_i F + 2P)/S + 1$ $D_{\alpha} = K$

Convolutional Neural Networks



VGG-16 Network



CNNs: Reminder Fully Connected Layers

Input

3072

(32 x 32 x 3 image -> stretches to 3072 x 1)





Convolutional Neural Networks



VGG-16 Network

CNNs: Reminder Fully Connected Layers



102,760,448 parameters!

* adopted from Fei-Dei Li, Justin Johnson, Serena Yeung, cs231n Stanford

Activation

4,096



Convolutional Neural Networks



VGG-16 Network

Pooling Layer



Let us assume the filter is an "eye" detector

How can we make detection spatially invariant (insensitive to position of the eye in the image)

* slide from Marc'Aurelio Renzato

Pooling Layer



Let us assume the filter is an "eye" detector

How can we make detection spatially invariant (insensitive to position of the eye in the image)

> By "pooling" (e.g., taking a max) response over a spatial locations we gain robustness to position variations



* slide from Marc'Aurelio Renzato

Pooling Layer

- Makes representation smaller, more manageable and spatially invariant
- Operates over each activation map independently



e manageable and spatially invariant independently



Max **Pooling**

activation map





max pool with 2 x 2 filter and stride of 2

6 8 3 4

Average **Pooling**

activation map





avg pool with 2 x 2 filter and stride of 2

3.25 5.25 2 2
Pooling Layer Receptive Field

If convolutional filters have size KxK and stride 1, and pooling layer has pools of size PxP, then each unit in the pooling layer depends upon a patch (at the input of the preceding conv. layer) of size: (P+K-1)x(P+K-1)



* slide from Marc'Aurelio Renzato

Pooling Layer Receptive Field

If convolutional filters have size KxK and stride 1, and pooling layer has pools of size PxP, then each unit in the pooling layer depends upon a patch (at the input of the preceding conv. layer) of size: **(P+K-1)x(P+K-1)**



* slide from Marc'Aurelio Renzato

Pooling Layer Summary

Accepts a volume of size: $W_i \times H_i \times D_i$ Requires hyperparameters: - Spatial extent of filters: K- Stride of application: FProduces a volume of size: $W_o \times H_o \times D_o$ $W_o = (W_i - F)/S + 1$ $H_o = (H_i - F)/S + 1$

Number of total learnable parameters: 0

$D_o = D_i$

Convolutional Neural Networks



VGG-16 Network

Local Contrast Normalization Layer

ensures response is the same in both case (details omitted, no longer popular)





* images from Marc'Aurelio Renzato