Lecture 7: Convolutional Networks

Justin Johnson

Lecture 7 - 1

Reminder: A2

Due Monday, September 30, 11:59pm (Even if you enrolled late!)

Your submission must pass the validation script

Justin Johnson

Lecture 7 - 2

Slight schedule change

Content originally planned for today got split into two lectures

Pushes the schedule back a bit:

A4 Due Date: Friday 11/1 -> Friday 11/8 A5 Due Date: Friday 11/15 -> Friday 11/22 A6 Due Date: Still Friday 12/6

Last Time: Backpropagation

Represent complex expressions as **computational graphs**



Forward pass computes outputs

Backward pass computes gradients

Lecture 7 - 4

During the backward pass, each node in the graph receives **upstream gradients** and multiplies them by **local gradients** to compute **downstream gradients**



September 24, 2019

Justin Johnson

f(x,W) = Wx





Input image (2, 2) **Problem**: So far our classifiers don't respect the spatial structure of images!

(4,)





Justin Johnson

Lecture 7 - 5

September 24, 2019

Stretch pixels into column

f(x,W) = Wx





Input image

(2, 2)

Problem: So far our classifiers don't respect the spatial structure of images!

56
231
24
2

(4,)





Solution: Define new computational nodes that operate on images!

Stretch pixels into column

Justin Johnson

Lecture 7 - 6

Components of a Full-Connected Network

Fully-Connected Layers



Activation Function



Justin Johnson

Lecture 7 - 7

Components of a Convolutional Network

Fully-Connected Layers



Activation Function



Convolution Layers



Pooling Layers



Normalization



Justin Johnson

Lecture 7 - 8

Components of a Convolutional Network

Fully-Connected Layers



Activation Function



Convolution Layers



Pooling Layers



Normalization



Justin Johnson

Lecture 7 - 9

Fully-Connected Layer

32x32x3 image -> stretch to 3072 x 1



Justin Johnson

Lecture 7 - 10

Fully-Connected Layer

32x32x3 image -> stretch to 3072 x 1



Justin Johnson

Lecture 7 - 11

Convolution Layer

3x32x32 image: preserve spatial structure



Justin Johnson

Lecture 7 - 12

Convolution Layer

3x32x32 image



3x5x5 filter

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

Justin Johnson

Lecture 7 - 13



Filters always extend the full depth of the input volume

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

Justin Johnson

Lecture 7 - 14

Convolution Layer

3x32x32 image



Justin Johnson

Lecture 7 - 15



Lecture 7 - 16



Lecture 7 - 17





Lecture 7 - 19



Lecture 7 - 20



Lecture 7 - 21



Lecture 7 - 22

Stacking Convolutions



Justin Johnson

Lecture 7 - 23

Stacking Convolutions

Q: What happens if we stack two convolution layers?



Justin Johnson

Lecture 7 - 24



Lecture 7 - 25



Justin Johnson

Lecture 7 - 26



Linear classifier: One template per class



Justin Johnson

Lecture 7 - 27



MLP: Bank of whole-image templates



Justin Johnson

Lecture 7 - 28



First-layer conv filters: local image templates (Often learns oriented edges, opposing colors)



AlexNet: 64 filters, each 3x11x11

Justin Johnson

Lecture 7 - 29



Justin Johnson

Lecture 7 - 30





7

Justin Johnson

Lecture 7 - 31



Input: 7x7 Filter: 3x3

7

Justin Johnson

Lecture 7 - 32





7

Justin Johnson

Lecture 7 - 33





7

Justin Johnson

Lecture 7 - 34



Input: 7x7 Filter: 3x3 Output: 5x5

7

Justin Johnson

Lecture 7 - 35



Input: 7x7 Filter: 3x3 Output: 5x5 **Problem: Feature** In general: maps "shrink" Input: W with each layer! Filter: K Output: W - K + 1

Justin Johnson

Lecture 7 - 36
A closer look at spatial dimensions

0	0	0	0	0	0	0	0	0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0	0	0	0	0	0	0	0	0

- Input: 7x7 Filter: 3x3 Output: 5x5
- In general:Problem: FeatureInput: Wmaps "shrink"Filter: Kwith each layer!

Output: W - K + 1

Solution: **padding** Add zeros around the input

Justin Johnson

Lecture 7 - 37

A closer look at spatial dimensions

0	0	0	0	0	0	0	0	0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0	0	0	0	0	0	0	0	0

Input: 7x7 Filter: 3x3 Output: 5x5

In general:Very common:Input: WSet P = (K - 1) / 2 toFilter: Kmake output havePadding: Psame size as input!

Output: W – K + 1 + 2P

Justin Johnson

For convolution with kernel size K, each element in the output depends on a K x K **receptive field** in the input



	U	S	ti	r		0	h	n	S	0	n
<u> </u>	<u> </u>		<u> </u>		• •				-	\smile	

Each successive convolution adds K - 1 to the receptive field size With L layers the receptive field size is 1 + L * (K - 1)



Input

Output

Be careful – "receptive field in the input" vs "receptive field in the previous layer" Hopefully clear from context!

```
Justin Johnson
```

Lecture 7 - 40

Each successive convolution adds K - 1 to the receptive field size With L layers the receptive field size is 1 + L * (K - 1)



Input

Problem: For large images we need many layers for each output to "see" the whole image image

Output

Justin Johnson

Lecture 7 - 41

Each successive convolution adds K - 1 to the receptive field size With L layers the receptive field size is 1 + L * (K - 1)



Input

Problem: For large images we need many layers for each output to "see" the whole image image Output

Solution: Downsample inside the network

Justin Johnson

Lecture 7 - 42

<u>Strided</u> Convolution



Input: 7x7 Filter: 3x3 Stride: 2

Justin Johnson

Lecture 7 - 43

Strided Convolution

Input: 7x7 Filter: 3x3 Stride: 2

Justin Johnson

Lecture 7 - 44

Strided Convolution

Input: 7x7 Filter: 3x3 Output: 3x3 Stride: 2

Justin Johnson

Lecture 7 - 45

<u>Strided</u> Convolution



Input: 7x7 Filter: 3x3 Output: 3x3 Stride: 2

In general: Input: W Filter: K Padding: P Stride: S Output: (W – K + 2P) / S + 1

Justin Johnson

Lecture 7 - 46

Input volume: 3 x 32 x 32 10 5x5 filters with stride 1, pad 2

Output volume size: ?



Justin Johnson

Lecture 7 - 47

Input volume: 3 x 32 x 32 10 5x5 filters with stride 1, pad 2

```
Output volume size:
(32+2*2-5)/1+1 = 32 spatially, so
10 x 32 x 32
```



Justin Johnson

Input volume: 3 x 32 x 32 10 5x5 filters with stride 1, pad 2

Output volume size: 10 x 32 x 32 Number of learnable parameters: ?



Justin Johnson

Lecture 7 - 49

Input volume: 3 x 32 x 32 10 5x5 filters with stride 1, pad 2

Output volume size: 10 x 32 x 32 Number of learnable parameters: **760** Parameters per filter: **3*5*5** + 1 (for bias) = **76 10** filters, so total is **10 * 76 = 760**



Input volume: 3 x 32 x 32 10 5x5 filters with stride 1, pad 2

Output volume size: 10 x 32 x 32 Number of learnable parameters: 760 Number of multiply-add operations: ?



Justin Johnson

Lecture 7 - 51

Input volume: **3** x 32 x 32 10 **5x5** filters with stride 1, pad 2



Output volume size: **10 x 32 x 32** Number of learnable parameters: 760 Number of multiply-add operations: **768,000 10*32*32** = 10,240 outputs; each output is the inner product of two **3**x**5**x**5** tensors (75 elems); total = 75*10240 = **768K**

Example: 1x1 Convolution



Justin Johnson

Lecture 7 - 53

Example: 1x1 Convolution



Convolution Summary

Input: C_{in} x H x W **Hyperparameters**:

- Kernel size: $K_H \times K_W$
- Number filters: C_{out}
- Padding: P
- Stride: S

Weight matrix: $C_{out} \times C_{in} \times K_H \times K_W$ giving C_{out} filters of size $C_{in} \times K_H \times K_W$ Bias vector: C_{out} Output size: $C_{out} \times H' \times W'$ where:

- H' = (H K + 2P) / S + 1
- W' = (W K + 2P) / S + 1

Convolution Summary

Input: C_{in} x H x W **Hyperparameters**:

- Kernel size: $K_H \times K_W$
- Number filters: C_{out}
- Padding: P
- Stride: S

Weight matrix: $C_{out} \times C_{in} \times K_H \times K_W$ giving C_{out} filters of size $C_{in} \times K_H \times K_W$ Bias vector: C_{out} Output size: $C_{out} \times H' \times W'$ where:

- H' = (H K + 2P) / S + 1
- W' = (W K + 2P) / S + 1

Common settings: $K_H = K_W$ (Small square filters) P = (K - 1) / 2 ("Same" padding) $C_{in}, C_{out} = 32, 64, 128, 256$ (powers of 2) K = 3, P = 1, S = 1 (3x3 conv) K = 5, P = 2, S = 1 (5x5 conv) K = 1, P = 0, S = 1 (1x1 conv) K = 3, P = 1, S = 2 (Downsample by 2)

Other types of convolution

So far: 2D Convolution



Justin Johnson

Lecture 7 - 57

Other types of convolution

So far: 2D Convolution



1D Convolution

Weights: C_{out} x C_{in} x K



Lecture 7 - 58

Other types of convolution

So far: 2D Convolution



C_{in}-dim vector at each point in the volume

3D Convolution

Input: C_{in} x H x W x D Weights: C_{out} x C_{in} x K x K x K



W

Justin Johnson

Lecture 7 - 59

PyTorch Convolution Layer

Conv2d

CLASS torch.nn.Conv2d(*in_channels*, *out_channels*, *kernel_size*, *stride=1*, *padding=0*, *dilation=1*, *groups=1*, *bias=True*, *padding_mode='zeros'*)

Applies a 2D convolution over an input signal composed of several input planes.

In the simplest case, the output value of the layer with input size $(N, C_{
m in}, H, W)$ and output $(N, C_{
m out}, H_{
m out}, W_{
m out})$ can be precisely described as:

$$\operatorname{out}(N_i, C_{\operatorname{out}_j}) = \operatorname{bias}(C_{\operatorname{out}_j}) + \sum_{k=0}^{C_{\operatorname{in}}-1} \operatorname{weight}(C_{\operatorname{out}_j}, k) \star \operatorname{input}(N_i, k)$$

Justin Johnson

Lecture 7 - 60

September 24, 2019

[SOURCE]

PyTorch Convolution Layers

Conv2d

Conv1d

[SOURCE]

[SOURCE]

Conv3d

[SOURCE]

Justin Johnson

Lecture 7 - 61

Components of a Convolutional Network

Fully-Connected Layers



Activation Function



Convolution Layers



Pooling Layers



Normalization



September 24, 2019

Justin Johnson

Lecture 7 - 62

Pooling Layers: Another way to downsample



Justin Johnson

Lecture 7 - 63

Justin Johnson

Х

Max Pooling

Single depth slice



Y

Max pooling with 2x2 kernel size and stride 2

Lecture 7 - 64



Introduces **invariance** to small spatial shifts No learnable parameters!



Pooling Summary

Input: C x H x W

Hyperparameters:

- Kernel size: K
- Stride: S
- Pooling function (max, avg)

Output: C x H' x W' where

- H' = (H K) / S + 1
- W' = (W K) / S + 1

Learnable parameters: None!

Common settings: max, K = 2, S = 2 max, K = 3, S = 2 (AlexNet)

Justin Johnson

Components of a Convolutional Network

Fully-Connected Layers



Activation Function



Convolution Layers



Pooling Layers



Normalization



Justin Johnson

Lecture 7 - 66

Convolutional Networks

Classic architecture: [Conv, ReLU, Pool] x N, flatten, [FC, ReLU] x N, FC



Lecun et al, "Gradient-based learning applied to document recognition", 1998

Justin Johnson

Lecture 7 - 67



Lecun et al, "Gradient-based learning applied to document recognition", 1998

Justin Johnson

Lecture 7 - 68



Layer	Output Size	Weight Size
Input	1 x 28 x 28	
Conv (C _{out} =20, K=5, P=2, S=1)	20 x 28 x 28	20 x 1 x 5 x 5
ReLU	20 x 28 x 28	



Lecun et al, "Gradient-based learning applied to document recognition", 1998

Justin Johnson

Lecture 7 - 69

Example: LeNet-5

Layer	Output Size	Weight Size
Input	1 x 28 x 28	
Conv (C _{out} =20, K=5, P=2, S=1)	20 x 28 x 28	20 x 1 x 5 x 5
ReLU	20 x 28 x 28	
MaxPool(K=2, S=2)	20 x 14 x 14	



Lecun et al, "Gradient-based learning applied to document recognition", 1998

Justin Johnson

Lecture 7 - 70

Example: LeNet-5

Layer	Output Size	Weight Size
Input	1 x 28 x 28	
Conv (C _{out} =20, K=5, P=2, S=1)	20 x 28 x 28	20 x 1 x 5 x 5
ReLU	20 x 28 x 28	
MaxPool(K=2, S=2)	20 x 14 x 14	
Conv (C _{out} =50, K=5, P=2, S=1)	50 x 14 x 14	50 x 20 x 5 x 5
ReLU	50 x 14 x 14	



Lecun et al, "Gradient-based learning applied to document recognition", 1998

Justin Johnson

Lecture 7 - 71

Example: LeNet-5

Layer	Output Size	Weight Size
Input	1 x 28 x 28	
Conv (C _{out} =20, K=5, P=2, S=1)	20 x 28 x 28	20 x 1 x 5 x 5
ReLU	20 x 28 x 28	
MaxPool(K=2, S=2)	20 x 14 x 14	
Conv (C _{out} =50, K=5, P=2, S=1)	50 x 14 x 14	50 x 20 x 5 x 5
ReLU	50 x 14 x 14	
MaxPool(K=2, S=2)	50 x 7 x 7	



Lecun et al, "Gradient-based learning applied to document recognition", 1998

Justin Johnson

Lecture 7 - 72
Layer	Output Size	Weight Size
Input	1 x 28 x 28	
Conv (C _{out} =20, K=5, P=2, S=1)	20 x 28 x 28	20 x 1 x 5 x 5
ReLU	20 x 28 x 28	
MaxPool(K=2, S=2)	20 x 14 x 14	
Conv (C _{out} =50, K=5, P=2, S=1)	50 x 14 x 14	50 x 20 x 5 x 5
ReLU	50 x 14 x 14	
MaxPool(K=2, S=2)	50 x 7 x 7	
Flatten	2450	



Lecun et al, "Gradient-based learning applied to document recognition", 1998

Justin Johnson

Lecture 7 - 73

Layer	Output Size	Weight Size
Input	1 x 28 x 28	
Conv (C _{out} =20, K=5, P=2, S=1)	20 x 28 x 28	20 x 1 x 5 x 5
ReLU	20 x 28 x 28	
MaxPool(K=2, S=2)	20 x 14 x 14	
Conv (C _{out} =50, K=5, P=2, S=1)	50 x 14 x 14	50 x 20 x 5 x 5
ReLU	50 x 14 x 14	
MaxPool(K=2, S=2)	50 x 7 x 7	
Flatten	2450	
Linear (2450 -> 500)	500	2450 x 500
ReLU	500	



Lecun et al, "Gradient-based learning applied to document recognition", 1998

Justin Johnson

Lecture 7 - 74

Layer	Output Size	Weight Size
Input	1 x 28 x 28	
Conv (C _{out} =20, K=5, P=2, S=1)	20 x 28 x 28	20 x 1 x 5 x 5
ReLU	20 x 28 x 28	
MaxPool(K=2, S=2)	20 x 14 x 14	
Conv (C _{out} =50, K=5, P=2, S=1)	50 x 14 x 14	50 x 20 x 5 x 5
ReLU	50 x 14 x 14	
MaxPool(K=2, S=2)	50 x 7 x 7	
Flatten	2450	
Linear (2450 -> 500)	500	2450 x 500
ReLU	500	
Linear (500 -> 10)	10	500 x 10





Lecture 7 - 75

September 24, 2019

Justin Johnson

Layer	Output Size	Weight Size
Input	1 x 28 x 28	
Conv (C _{out} =20, K=5, P=2, S=1)	20 x 28 x 28	20 x 1 x 5 x 5
ReLU	20 x 28 x 28	
MaxPool(K=2, S=2)	20 x 14 x 14	
Conv (C _{out} =50, K=5, P=2, S=1)	50 x 14 x 14	50 x 20 x 5 x 5
ReLU	50 x 14 x 14	
MaxPool(K=2, S=2)	50 x 7 x 7	
Flatten	2450	
Linear (2450 -> 500)	500	2450 x 500
ReLU	500	
Linear (500 -> 10)	10	500 x 10



As we go through the network:

Spatial size **decreases** (using pooling or strided conv)

Number of channels **increases** (total "volume" is preserved!)

Lecun et al, "Gradient-based learning applied to document recognition", 1998

Justin Johnson

Problem: Deep Networks very hard to train!

Justin Johnson

Lecture 7 - 77

Components of a Convolutional Network

Fully-Connected Layers



Activation Function



Convolution Layers



Pooling Layers



Normalization



Justin Johnson

Lecture 7 - 78

Idea: "Normalize" the outputs of a layer so they have zero mean and unit variance

Why? Helps reduce "internal covariate shift", improves optimization

We can normalize a batch of activations like this:



This is a **differentiable function**, so we can use it as an operator in our networks and backprop through it!

loffe and Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift", ICML 2015

Justin Johnson



loffe and Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift", ICML 2015

Justin Johnson

Lecture 7 - 80

Input: $x: N \times D$



 $\mu_j = \frac{1}{N} \sum_{i=1}^{N} x_{i,j}$ Per-channel mean, shape is D $\sigma_j^2 = \frac{1}{N} \sum_{i=1}^{N} (x_{i,j} - \mu_j)^2 \quad \begin{array}{l} \text{Per-channel} \\ \text{std, shape is D} \end{array}$ i=1 $\hat{x}_{i,j} = rac{x_{i,j} - \mu_j}{\sqrt{\sigma_i^2 + \varepsilon}}$ Normalized x, Shape is N x D

)

Problem: What if zero-mean, unit variance is too hard of a constraint?

loffe and Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift", ICML 2015

Justin Johnson

Ν

Lecture 7 - 81

Input: $x: N \times D$

Learnable scale and shift parameters:

 $\gamma,\beta:D$

Learning $\gamma = \sigma$, $\beta = \mu$ will recover the identity function!

$$\begin{split} \mu_j &= \frac{1}{N} \sum_{i=1}^N x_{i,j} & \text{Per-channel} \\ \text{mean, shape is D} \\ \sigma_j^2 &= \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2 & \text{Per-channel} \\ \text{std, shape is D} \\ \hat{x}_{i,j} &= \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}} & \text{Normalized x,} \\ y_{i,j} &= \gamma_j \hat{x}_{i,j} + \beta_j & \text{Output,} \end{split}$$

A T

Shape is N x D

Problem: Estimates depend on Batch Normalization: Test-Time minibatch; can't do this at test-time!

Input: $x: N \times D$

Learnable scale and shift parameters:

 $\gamma, \beta: D$

Learning $\gamma = \sigma$, $\beta = \mu_{\rm e}$ will recover the identity function!

$$\mu_{j} = \frac{1}{N} \sum_{i=1}^{N} x_{i,j} \quad \begin{array}{l} \text{Per-channel} \\ \text{mean, shape is D} \end{array}$$

$$\sigma_{j}^{2} = \frac{1}{N} \sum_{i=1}^{N} (x_{i,j} - \mu_{j})^{2} \quad \begin{array}{l} \text{Per-channel} \\ \text{std, shape is D} \end{array}$$

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_{j}}{\sqrt{\sigma_{i}^{2} + \varepsilon}} \quad \begin{array}{l} \text{Normalized x,} \\ \text{Shape is N x D} \end{array}$$

$$rac{p-\mu_j}{\sigma_j^2+arepsilon}$$
 Normaliz

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$
 Out

itput, Shape is N x D

N x D

Batch Normalization: Test-Time

Input: $x: N \times D$

Learnable scale and shift parameters:

 $\gamma, \beta: D$

Learning $\gamma = \sigma$, $\beta = \mu$ will recover the identity function!

(Running) average of Per-channel $\mu_j = \text{values seen during}$ mean, shape is D training (Running) average of $\sigma_j^2 =$ values seen during Per-channel std, shape is D training $\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$ Normalized x, Shape is N x D $y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$ Output, Shape is N x D

Batch Normalization: Test-Time

Input: $x: N \times D$

Learnable scale and shift parameters:

 $\gamma,\beta:D$

During testing batchnorm becomes a linear operator! Can be fused with the previous fully-connected or conv layer (Running) average of $\mu_j = ext{values seen during}$ training

 $\sigma_j^2 = \begin{array}{l} \text{(Running) average of} \\ \text{values seen during} \\ \text{training} \end{array}$

Per-channel mean, shape is D

Per-channel std, shape is D

 $\hat{x}_{i,j} = rac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$ No

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$

Normalized x, Shape is N x D

Output, Shape is N x D

Batch Normalization for ConvNets

Batch Normalization for **fully-connected** networks

Batch Normalization for **convolutional** networks (Spatial Batchnorm, BatchNorm2D)

X :	Ν	X	D	
Normalize	Ļ			
μ,σ:	1	X	D	
γ,β:	1	X	D	
y = y	x (>	<u>د – ا</u>	u)/(7 +β

x: N×C×H×W Normalize \downarrow \downarrow \downarrow μ, σ : 1×C×1×1 γ, β : 1×C×1×1 $y = \gamma(x-\mu)/\sigma+\beta$



Usually inserted after Fully Connected or Convolutional layers, and before nonlinearity.



loffe and Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift", ICML 2015

Justin Johnson

Lecture 7 - 87



- Makes deep networks **much** easier to train!
- Allows higher learning rates, faster convergence
- Networks become more robust to initialization
- Acts as regularization during training
- Zero overhead at test-time: can be fused with conv!



loffe and Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift", ICML 2015

Justin Johnson

Lecture 7 - 88

_

-



- Makes deep networks **much** easier to train!
- Allows higher learning rates, faster convergence
- Networks become more robust to initialization
- Acts as regularization during training
- Zero overhead at test-time: can be fused with conv!
- Not well-understood theoretically (yet)
- Behaves differently during training and testing: this is a very common source of bugs!

loffe and Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift", ICML 2015

Justin Johnson

Lecture 7 - 89

Layer Normalization

Batch Normalization for fully-connected networks



Layer Normalization for fullyconnected networks Same behavior at train and test! Used in RNNs, Transformers

X :	Ν	×	D		
Normalize			Ţ		
μ,σ:	Ν	X	1		
γ,β:	1	X	D		
y = y	x (>	د – ا	4) /	/σ+ β	3

Ba, Kiros, and Hinton, "Layer Normalization", arXiv 2016

Justin Johnson

Instance Normalization

Batch Normalization for convolutional networks

Instance Normalization for convolutional networks Same behavior at train / test!



Ulyanov et al, Improved Texture Networks: Maximizing Quality and Diversity in Feed-forward Stylization and Texture Synthesis, CVPR 2017

Justin Johnson

Lecture 7 - 91

Comparison of Normalization Layers



Wu and He, "Group Normalization", ECCV 2018

Justin Johnson

Lecture 7 - 92

Group Normalization



Wu and He, "Group Normalization", ECCV 2018

Justin Johnson

Lecture 7 - 93

Components of a Convolutional Network

Convolution Layers



Pooling Layers



Fully-Connected Layers



Activation Function



Normalization

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

September 24, 2019

Justin Johnson

Components of a Convolutional Network



Pooling Layers



Fully-Connected Layers



Activation Function



Normalization

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

September 24, 2019

Justin Johnson

Summary: Components of a Convolutional Network

Convolution Layers



Pooling Layers



Fully-Connected Layers



Activation Function



Normalization

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

September 24, 2019

Justin Johnson

Summary: Components of a Convolutional Network

Problem: What is the right way to combine all these components?



		S.	ti	n	lo	h	n	S	\cap	n	
J	U	5	U					9	$\mathbf{\circ}$		

Next time: CNN Architectures

Justin Johnson

Lecture 7 - 98