Topics:

Convolutional Neural Networks

CS 4644-DL / 7643-A ZSOLT KIRA

Assignment 2 – Due June 22nd

- Implement convolutional neural networks
- Resources (in addition to lectures):
 - DL book: Convolutional Networks
 - CNN notes https://www.cc.gatech.edu/classes/AY2022/cs7643 spring/assets/L10 cnns notes.pdf
 - Backprop notes https://www.cc.gatech.edu/classes/AY2023/cs7643_spring/assets/L10_cnns_backprop_notes.pdf
 - HW2 Tutorial, Conv backward Coming Soon
 - Slower OMSCS lectures on dropbox: Module 2 Lessons 5-6 (M2L5/M2L6) (https://www.dropbox.com/sh/iviro188gq0b4vs/AADdHxX_Uy1TkpF_yvIzX0nPa?dl=0)
- Project
 - Proposal planning session Wed.
 - Proposal due June 15th
- **Reminder:** Please do readings announced for discussions!

$$X(0:2,0:2) = \begin{bmatrix} 200 \ 150 \ 100 \\ 100 \ 50 \ 100 \\ 25 \ 25 \ 10 \end{bmatrix} \qquad K' = \begin{bmatrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{bmatrix} \xrightarrow{\text{Dot product}} X(0:2,0:2) \cdot K' = 65 + \text{bias}$$
Dot product (element-wise multiply and sum)
$$y(r,c) = (x * k)(r,c) = \sum_{a=0}^{k_1-1} \sum_{b=0}^{k_2-1} x(r + a, c + b) k(a, b)$$

Cross-Correlation



Number of parameters with N filters is: $N * (k_1 * k_2 * 3 + 1)$

Number of Parameters





Need to incorporate all upstream gradients:

 $\left\{\frac{\partial L}{\partial y(0,0)}, \frac{\partial L}{\partial y(0,1)}, \dots, \frac{\partial L}{\partial y(H,W)}\right\}$

Chain Rule: $\frac{\partial L}{\partial k(a,b)} = \sum_{r=0}^{H-1} \sum_{c=0}^{W-1} \frac{\partial L}{\partial y(r,c)} \frac{\partial y(r,c)}{\partial k(a,b)}$ Sum over Upstream We will all output gradient compute pixels (known)





 $W=5 \qquad (H-1,W-1)$





 $\frac{\partial y(r,c)}{\partial k(a,b)} = x(r+a,c+b)$

$$\frac{\partial L}{\partial k(a,b)} = \sum_{r=0}^{H-1} \sum_{c=0}^{W-1} \frac{\partial L}{\partial y(r,c)} x(r+a,c+b)$$



Gradients and Cross-Correlation

Georgia Tech $\frac{\partial y(r,c)}{\partial k(a,b)} = x(r+a,c+b)$

$$\frac{\partial L}{\partial k(a,b)} = \sum_{r=0}^{H-1} \sum_{c=0}^{W-1} \frac{\partial L}{\partial y(r,c)} x(r+a,c+b)$$

Does this look familiar?

Cross-correlation between upstream gradient and input! (until $k_1 \times k_2$ output)





Gradients and Cross-Correlation

Georg a Tech $\frac{\partial L}{\partial x} = \frac{\partial L}{\partial y} \quad \frac{\partial y}{\partial x}$

Gradient for input (to pass to prior layer)

Calculate one pixel at a time

$$\frac{\partial L}{\partial x(r',c')}$$

What does this input pixel affect at the output?

Neighborhood around it (where part of the kernel touches it)





 $W = 5 \qquad (H-1, W-1)$



What an Input Pixel Affects at Output

Chain rule for affected pixels (sum gradients):





Summing Gradient Contributions



Chain rule for affected pixels (sum gradients):

Let's derive it analytically this time (as opposed to visually)





Summing Gradient Contributions

Plugging in to earlier equation:

$$\frac{\partial L}{\partial x(r',c')} = \sum_{a=0}^{k_1-1} \sum_{b=0}^{k_2-1} \frac{\partial L}{\partial y(r'-a,c'-b)} \frac{\partial y(r'-a,c'-b)}{\partial x(r',c')}$$

Does this look familiar?

$$=\sum_{a=0}^{k_1-1}\sum_{b=0}^{k_2-1}\frac{\partial L}{\partial y(r'-a,c'-b)}k(a,b)$$

Again, all operations can be implemented via matrix multiplications (same as FC layer)! Convolution between upstream gradient and kernel!

(can implement by flipping kernel and cross- correlation)





- Convolutions are mathematical descriptions of striding linear operation
- In practice, we implement **cross-correlation neural networks!** (still called convolutional neural networks due to history)
 - Can connect to convolutions via duality (flipping kernel)
 - Convolution formulation has mathematical properties explored in ECE
- Duality for forwards and backwards:
 - Forward: Cross-correlation
 - Backwards w.r.t. K: Cross-correlation b/w upstream gradient and input
 - Backwards w.r.t. X: Convolution b/w upstream gradient and kernel
 - In practice implement via cross-correlation and flipped kernel
- All operations still implemented via **efficient linear algebra** (e.g. matrixmatrix multiplication)





Pooling Layers



- Dimensionality reduction is an important aspect of machine learning
- Can we make a layer to explicitly down-sample image or feature maps?



Yes! We call one class of these operations pooling operations

Parameters

- kernel_size the size of the window to take a max over
- stride the stride of the window. Default value is kernel_size
- padding implicit zero padding to be added on both sides

From: https://pytorch.org/docs/stable/generated/torch.nn.MaxPool2d.html#torch.nn.MaxPool2/





Example: Max pooling

Stride window across image but perform per-patch max operation

 $X(0:2,0:2) = \begin{bmatrix} 200 & 150 & 150 \\ 100 & 50 & 100 \\ 25 & 25 & 10 \end{bmatrix} \implies \max(0:2,0:2) = 200$



Max Pooling

How many learned parameters does this layer have?





- makes the representations spatially smaller
- saves computation (GPU mem & speed), allows go deeper
- operates over each activation map independently:



From: Slides by CS 231n, Danfei Xu



Since the **output** of convolution and pooling layers are **(multi-channel) images**, we can sequence them just as any other layer







This combination adds some invariance to translation of the features

If feature (such as beak) translated a little bit, output values still remain the same







Convolution by itself has the property of equivariance

If feature (such as beak) translated a little bit, output values move by the same translation





W = 5



Invariance vs. Equivariance

Simple Convolutional Neural Networks



Since the **output** of convolution and pooling layers are **(multi-channel) images**, we can sequence them just as any other layer



Combining Convolution & Pooling Layers



Alternating Convolution and Pooling





Geo







Receptive Fields



These architectures have existed **since 1980s**



Image Credit: Yann LeCun, Kevin Murphy



Georg a Tech

Handwriting Recognition



Image Credit: Yann LeCun Georg a

Translation Equivariance (Conv Layers) & Invariance (Output)



Image Credit: Yann LeCun Georgia

(Some) Rotation Invariance



Image Credit: Yann LeCun Georga

(Some) Scale Invariance



Image Credit: Yann LeCun Georgaa







Advanced Convolutional Networks





The **ImageNet** dataset contains 14,197,122 annotated images according to the WordNet hierarchy. ImageNet Large Scale Visual Recognition Challenge (ILSVRC) is a benchmark for image classification and object detection based on the dataset.

Benchmarking Models

The Importance of Benchmarks



C



Case Studies

- AlexNet
- VGG
- GoogLeNet
- ResNet



Also....

- SENet
- Wide ResNet
- ResNeXT

- DenseNet
- MobileNets
- NASNet
- EfficientNet
- ConvNeXt v1/v2

The Space of CNN Architectures



AlexNet - Architecture



From: Krizhevsky et al., ImageNet Classification with Deep ConvolutionalNeural Networks, 2012.


Case Study: AlexNet

[Krizhevsky et al. 2012]





Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.







First layer (CONV1): 96 11x11 filters applied at stride 4 => Q: what is the output volume size? Hint: (227-11)/4+1 = 55 W' = (W - F + 2P) / S + 1

From: Slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231r





Output volume [55x55x96]

First layer (CONV1): 96 11x11 filters applied at stride 4

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

227 227 3 55×55 96

From: Slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n





```
First layer (CONV1): 96 11x11 filters applied at stride 4
=>
Output volume [55x55x96]
```

Q: What is the total number of parameters in this layer?



From: Slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231r





Parameters: (11*11*3 + 1)*96 = 35K

```
First layer (CONV1): 96 11x11 filters applied at stride 4
=>
Output volume [55x55x96]
```

11 x 11

From: Slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231r



Full (simplified) AlexNet architecture: [224k224x3] INPUT [55x55x96] CONV1: 96 11x11 filters at stride 4, pad 0 [27x27x96] MAX POOL1: 3x3 filters at stride 2 [27x27x96] NORM1: Normalization layer [27x27x256] CONV2: 256 5x5 filters at stride 1, pad 2 [13x13x256] MAX POOL2: 3x3 filters at stride 2 [13x13x256] NORM2: Normalization layer [13x13x384] CONV3: 384 3x3 filters at stride 1, pad 1 [13x13x384] CONV4: 384 3x3 filters at stride 1, pad 1 [13x13x256] MAX POOL3: 3x3 filters at stride 1, pad 1 [13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1 [6x6x256] MAX POOL3: 3x3 filters at stride 2 [4096] FC6: 4096 neurons [4096] FC7: 4096 neurons [1000] FC8: 1000 neurons (class scores)



Key aspects:

- ReLU instead of sigmoid or tanh
- Specialized normalization layers
- PCA-based data augmentation
- Dropout
- Ensembling

From: Slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231r

Georgia Tech

Small filters, Deeper networks

8 layers (AlexNet) -> 16 - 19 layers (VGG16Net)

Only 3x3 CONV stride 1, pad 1 and 2x2 MAX POOL stride 2

11.7% top 5 error in ILSVRC'13 (ZFNet) -> 7.3% top 5 error in ILSVRC'14









Q: Why use smaller filters? (3x3 conv)

Stack of three 3x3 conv (stride 1) layers has same effective receptive field as one 7x7 conv layer

Q: What is the effective receptive field of three 3x3 conv (stride 1) layers?



FC 4096

C 4096





Q: What is the effective receptive field of three 3x3 conv (stride 1) layers?



Conv1 (3x3) Conv2 (3x3) Conv3 (3x3)



But deeper, more non-linearities

And fewer parameters: $3 * (3^2C^2)$ vs. 7^2C^2 for C channels per layer





INDUT: (224x224x21 memory: 224x224x2=150K percent; 0 (not counting biases)	ConvNet Configuration						
inPO1. [224x224x3] memory. 224 224 3=150K params. 0 (net counting biddece)		А	A-LRN	В	С	D	Е
CONV3-64: [224x224x64] memory: 224*224*64=3.2M params: (3*3*3)*64 = 1,728	11	weight	11 weight	13 weight	16 weight	16 weight	19 weight
CONV3-64: [224x224x64] memory: 224*224*64=3.2M params: (3*3*64)*64 = 36,864		ayers	layers	layers	layers	layers	layers
POOL2 [112x112x64] memory 112*112*64=800K params 0			input (224×224 RGB image)				
CONV3-128: [112x112x128], memory: 112*112*128=1.6M, params: (3*3*64)*128 = 73.728	co	nv3-64	conv3-64	conv3-64	conv3-64	conv3-64	conv3-64
CONV2 120: [112x112x120] memory 1121121120=1.0M parameter (21214)20120 = 10,720			LRN	conv3-64	conv3-64	conv3-64	conv3-64
CONV3-128 [112x112x128] memory. 112 112 128 1.6M params. (3 3 128) 128 = 147,456		2 100		max	pool	100	120
POOL2: [56x56x128] memory: 56*56*128=400K params: 0	con	IV3-128	conv3-128	conv3-128	conv3-128	conv3-128	conv3-128
CONV3-256: [56x56x256] memory: 56*56*256=800K params: (3*3*128)*256 = 294,912				tonvo-128	pool	conv3-128	conv5-128
CONV3-256: [56x56x256] memory: 56*56*256=800K params: (3*3*256)*256 = 589.824	COL	w3-256	conv3-256	conv3-256	conv3-256	conv3-256	conv3-256
CONV/3-256: [56x56x256], memory: 56*56*256=800K, params: (3*3*256)*256 = 589.824	con	w3-256	conv3-256	conv3-256	conv3-256	conv3-256	conv3-256
DOOL 2: [29/2252] memory: 29/29/256 = 200K prome: 0					conv1-256	conv3-256	conv3-256
POULZ. [282282286] Memory. 28 28 286-200K params. 0							conv3-256
CONV3-512: $[28x28x512]$ memory: $28^{\circ}28^{\circ}512 = 400K$ params: $(3^{\circ}3^{\circ}256)^{\circ}512 = 1,179,648$				max	pool		
CONV3-512: [28x28x512] memory: 28*28*512=400K params: (3*3*512)*512 = 2,359,296	COL	w3-512	conv3-512	conv3-512	conv3-512	conv3-512	conv3-512
CONV3-512: [28x28x512] memory: 28*28*512=400K params: (3*3*512)*512 = 2.359.296	COL	w3-512	conv3-512	conv3-512	conv3-512	conv3-512	conv3-512
POOL 2: [14x14x512] memory: 14*14*512=100K params: 0					conv1-512	conv5-512	conv3 512
CON(2, 512) [14/14/14/512] memory: 14/14/512=100K, parame: (2*2*512)*512 = 2.250.206	e maxpool			pool		conv3-312	
CONV3-512. [14114/312] memory. 14 14 512-100K params. (3 3 512) 512 - 2,359,290	COL	w3-512	conv3-512	conv3-512	conv3-512	conv3-512	conv3-512
CONV3-512: $[14x14x512]$ memory: $14*14*512=100K$ params: $(3*3*512)*512 = 2,359,296$	con	w3-512	conv3-512	conv3-512	conv3-512	conv3-512	conv3-512
CONV3-512: [14x14x512] memory: 14*14*512=100K params: (3*3*512)*512 = 2,359,296					conv1-512	conv3-512	conv3-512
POOL2: [7x7x512] memory: 7*7*512=25K params: 0							conv3-512
EC: $[1x1x4096]$ memory: 4096 params: $7*7*512*4096 = 102760448$	maxpool						
C_{1} [1x1x4000] memory 4000 parame: 4000*4000 = 16 777.216	FC-4096						
PC. [1x1x4030] memory. 4036 params. 4036 4036 - 16,777,216	FC-4096						
FC: $[1x1x1000]$ memory: 1000 params: 4096*1000 = 4,096,000	FC-1000						
	son-max						

VGG

Table 2: Number of parameters (in millions).							
Network	A,A-LRN	В	С	D	E		
Number of parameters	133	133	134	138	144		

From: Simonyan & Zimmerman, Very Deep Convolutional Networks for Large-Scale Image Recognition From: Slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231r/



INPUT: [224x224x3] memory: 224*224*3=150K params: 0 (not counting biases)							
CONV3-64: [224x224x64] memory: 224*224*64=3.2M params: (3*3*3)*64 = 1,728							
CONV3-64: [224x224x64] memory: 224*224*64=3.2M params: (3*3*64)*64 = 36,864							
POOL2: [112x112x64] memory: 112*112*64=800K params: 0							
CONV3-128: [112x112x128] memory: 112*112*128=1.6M params: (3*3*64)*128 = 73,728							
CONV3-128: [112x112x128] memory: 112*112*128=1.6M params: (3*3*128)*128 = 147,456							
POOL2: [56x56x128] memory: 56*56*128=400K params: 0							
CONV3-256: [56x56x256] memory: 56*56*256=800K params: (3*3*128)*256 = 294,912							
CONV3-256: [56x56x256] memory: 56*56*256=800K params: (3*3*256)*256 = 589,824							
CONV3-256: [56x56x256] memory: 56*56*256=800K params: (3*3*256)*256 = 589,824							
POOL2: [28x28x256] memory: 28*28*256=200K params: 0							
CONV3-512: [28x28x512] memory: 28*28*512=400K params: (3*3*256)*512 = 1,179,648							
CONV3-512: [28x28x512] memory: 28*28*512=400K params: (3*3*512)*512 = 2,359,296							
CONV3-512: [28x28x512] memory: 28*28*512=400K params: (3*3*512)*512 = 2,359,296							
POOL2: [14x14x512] memory: 14*14*512=100K params: 0							
CONV3-512: [14x14x512] memory: 14*14*512=100K params: (3*3*512)*512 = 2,359,296							
CONV3-512: [14x14x512] memory: 14*14*512=100K params: (3*3*512)*512 = 2,359,296							
CONV3-512: [14x14x512] memory: 14-14-512=100K params: (3-3-512)-512 = 2,359,296							
POOL2. [7x7x512] memory: 777512=25K params: 0							
FC. $[1x1x4096]$ memory: 4096 parame: 4096*4096 = 102,700,446							
FC: [1x1x4030] memory: 1000 parame: $4030 = 10,777,210$							
-6. [1x1x1000] memory. 1000 params. 4090 1000 = 4,090,000							

Most memory usage in convolution layers

Most parameters in FC layers

From: Simonyan & Zimmerman, Very Deep Convolutional Networks for Large-Scale Image Recognition From: Slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231r



Parameters and Memory

Case Study: VGGNet

[Simonyan and Zisserman, 2014]

Details:

- ILSVRC'14 2nd in classification, 1st in localization
- Similar training procedure as Krizhevsky 2012
- No Local Response Normalisation (LRN)
- Use VGG16 or VGG19 (VGG19 only slightly better, more memory)
- Use ensembles for best results
- FC7 features generalize well to other tasks

Still very expensive!

TOTAL memory: 24M * 4 bytes ~= 96MB / image (only forward! ~*2 for bwd) TOTAL params: 138M parameters

VGGNet

			Softmax
			FC 1000
		Softmax	FC 4098
	fc8	FC 1000	FC 4098
	fc7	FC 4098	Pool
	fc6	FC 4098	3x3 conv, 512
		Pool	3x3 conv. 512
	conv5-3	3x3 conv, 512	3x3 conv. 512
	conv5-2	3x3 conv, 512	3x3 conv. 512
	conv5-1	3x3 conv. 512	Pool
		Pool	3x3 conv. 512
Softmax	conv4-3	3x3 conv. 512	3x3 conv, 512
FC 1000	conv4-2	3x3 conv, 512	3x3 conv, 512
FC 4098	conv4-1	3x3 conv, 512	3x3 conv. 512
FC 4098		Pool	Pool
Pool	conv3-2	3x3 conv, 258	3x3 conv, 256
3x3 conv. 256	conv3-1	3x3 conv. 258	3x3 conv. 256
3x3 conv. 384		Pool	Pool
Pool	conv2-2	3x3 conv. 128	3x3 conv. 128
3x3 conv. 384	conv2-1	3x3 conv. 128	3x3 conv. 128
Pool		Pool	Pool
5x5 conv, 258	conv1-2	3x3 conv, 64	3x3 conv, 64
1x11 conv, 96	conv1-1	3x3 conv, 64	3x3 conv, 64
Input		Input	Input
AlexNet		VGG16	VGG19



fc7

fc6

conv5

conv4

conv3

conv2 conv1



But have become **deeper and more complex**



From: Szegedy et al. Going deeper with convolutions



Georg a Tech

Case Study: GoogLeNet

[Szegedy et al., 2014]

Deeper networks, focus on computational efficiency

- ILSVRC'14 classification winner (6.7% top 5 error)
- 22 layers
- Only 5 million parameters!
 12x less than AlexNet
 27x less than VGG-16
- Efficient "Inception" module
- No FC layers



From: Szegedy et al. Going deeper with convolutions



Case Study: GoogLeNet

[Szegedy et al., 2014]

Deeper networks, focus on computational efficiency

- ILSVRC'14 classification winner (6.7% top 5 error)
- 22 layers
- Only 5 million parameters!
 12x less than AlexNet
 27x less than VGG-16
- Efficient "Inception" module
- No FC layers





From: Szegedy et al. Going deeper with convolutions



Inception Architecture

Key idea: Repeated blocks and multi-scale features



From: Szegedy et al. Going deeper with convolutions





Case Study: GoogLeNet

[Szegedy et al., 2014]

"Inception module": design a good local network topology (network within a network) and then stack these modules on top of each other

Multiple conv filter size diversifies learned features



















Alternatively, interpret it as applying the same FC layer on each input pixel









1x1 Convolutions

Case Study: GoogLeNet

[Szegedy et al., 2014]







Case Study: GoogLeNet

[Szegedy et al., 2014]

Deeper networks, with computational efficiency

- 22 layers
- Efficient "Inception" module
- Avoids expensive FC layers
- 12x less params than AlexNet
- 27x less params than VGG-16
- ILSVRC'14 classification winner (6.7% top 5 error)

Fiter concatenation Previous Lave Inception module

From: slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n



69











The Challenge of Depth



From: He et al., Deep Residual Learning for Image Recognition

Optimizing very deep networks is challenging!



[He et al., 2015]

A deeper model can **emulate** a shallower model: copy layers from shallower model, set extra layers to identity

Thus deeper models should do at least <u>as good as</u> shallow models

Deeper models are harder to optimize. They don't learn identity functions (no-op) to emulate shallow models

Solution: Change the network so learning identity functions (no-op) as extra layers is easy







[He et al., 2015]

Solution: Change the network so learning identity functions as extra layers is easy







Case Study: ResNet

[He et al., 2015]

Full ResNet architecture:

- Stack residual blocks
- Every residual block has two 3x3 conv layers



FC 1000





Case Study: ResNet

[He et al., 2015]

Full ResNet architecture:

- Stack residual blocks
- Every residual block has two 3x3 conv layers
- Periodically, double # of filters and downsample spatially using stride 2 (/2 in each dimension) Reduce the activation volume by half.



oftmax

FC 1000





Case Study: ResNet

[He et al., 2015]

Full ResNet architecture:

- Stack residual blocks
- Every residual block has two 3x3 conv layers
- Periodically, double # of filters and downsample spatially using stride 2 (/2 in each dimension) Reduce the activation volume by half.
- Additional conv layer at the beginning (stem)



FC 1000

Skip Conections



[He et al., 2015]

For deeper networks (ResNet-50+), use "bottleneck" layer to improve efficiency (similar to GoogLeNet)







				<pre>>>> import torch >>> from torchvision. >>> model = resnet18(>>> summary(model2, (</pre>	models import resnet18) 3, 224, 224), device='d	cpu')
layer name	output size	18-layer	34-layer	Layer (type)	Output Sł	nape Param
conv1	112×112			conv2d 1	 Γ 1 64 112 1	
				BatchNorm2d-2	[-1, 64, 112, 1 [-1, 64, 112, 1	[12] 9,40 [12] 12
				Rel U-3	[-1, 64, 112, 1]	112]
$conv2_x$ 56×56	56×56	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 2$	$\left[\begin{array}{c} 3\times3,64\\ 3\times3,64 \end{array}\right]\times3$	MaxPool2d-4	[-1, 64, 56,	56]
	30730			Conv2d-5	[-1, 64, 56,	56] 36,86
				BatchNorm2d-6	[-1, 64, 56,	56] 12
			ReLU-7	[-1, 64, 56,	56]	
$conv3$ x 28×28	$\begin{bmatrix} 3 \times 3, 128 \end{bmatrix}_{\times 2}$	[3×3, 128] ×4	Conv2d-8	[-1, 64, 56,	56] 36,86	
			BatchNorm2d-9	[-1, 64, 56,	56] 12	
convo_x	20720	3×3, 128 ^2	3×3, 128 ^	ReLU-10	[-1, 64, 56,	56]
				BasicBlock-11	[-1, 64, 56,	56]
			$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 6$	Conv2d-12	[-1, 64, 56,	56] 36,86
	14 14	$4 \times 14 \begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 2$		BatchNorm2d-13	[-1, 64, 56,	56] 12
conv4_x	14×14			ReLU-14	[-1, 64, 56,	56]
			Conv2d-15	[-1, 64, 56,	56 36,86	
				BatchNorm2d-16	[-1, 64, 56,	56] 12
$conv5_x$ 7×7	[3×3 512]	[3×3 512]	ReLU-17 RecicPlack 19	[-1, 04, 30, [1 64 F6	50] 56]	
	7×7	$ \frac{5 \times 5, 512}{2} \times 2$	$ 3 \times 3, 512 \times 3$		[-1, 04, 30, [_1 130 30	טכ 201 ד ד סכ
		[3×3, 512]	[3×3, 512]	BatchNorm2d-20	$\begin{bmatrix} -1 & 120 & 20 \end{bmatrix}$	28] 73,72
				Rel U-21	[-1, 128, 28]	28]
	1×1		ave	Conv2d-22	[-1, 128, 28,	28] 147,45
FLO	OPs	1.8×10^{9}	3.6×10^{9}	3.8×10°	7.6×10 ²¹²⁰ 20	11.3×10°25





Training ResNet in practice:

- Batch Normalization after every CONV layer
- Xavier initialization from He et al.
- SGD + Momentum
- Learning rate: 0.1, divided by 10 when validation error plateaus
- Mini-batch size 256
- Weight decay of 1e-5
- No dropout used





Computational Complexity





0

GE

From: An Analysis Of Deep Neural Network Models For Practical Application

Wide Residual Networks

[Zagoruyko et al. 2016]

- Argues that residuals are the important factor, not depth
- Use wider residual blocks (F x k filters instead of F filters in each layer)
- 50-layer wide ResNet outperforms
 152-layer original ResNet
- Increasing width instead of depth more computationally efficient (parallelizable)







Densely Connected Convolutional Networks (DenseNet)

[Huang et al. 2017]

- Dense blocks where each layer is connected to every other layer in feedforward fashion
- Alleviates vanishing gradient, strengthens feature propagation, encourages feature reuse
- Showed that shallow 50-layer network can outperform deeper 152 layer ResNet

DenseNet



Dense Block


ConvNeXt (2022)

• Do the reading for discussion!



Slides created for CS886 at UWaterloo



Several ways to *learn* architectures:

- Evolutionary learning and reinforcement learning
- Prune overparameterized networks
- Learning of repeated blocks typical



From: https://ai.googleblog.com/2018/03/using-evolutionary-automl-to-discover.html





Learning Transferable Architectures for Scalable Image Recognition

[Zoph et al. 2017]

- Applying neural architecture search (NAS) to a large dataset like ImageNet is expensive
- Design a search space of building blocks ("cells") that can be flexibly stacked
- NASNet: Use NAS to find best cell structure on smaller CIFAR-10 dataset, then transfer architecture to ImageNet









2 x 2 maxpoo

Evolving Architectures and AutoML

 Convolutional neural networks (CNNs) stack pooling, convolution, nonlinearities, and fully connected (FC) layers

Feature engineering => architecture engineering!

- Tons of small details and tips/tricks
- Considerations: Memory, compute/FLO, dimensionality reduction, diversity of features, number of parameters/capacity, etc.





Transfer Learning & Generalization





From: slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n







From: slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n



Georgia Tech



From: slides by Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n





What if we don't have enough data?

Step 1: Train on large-scale dataset







Networks

Transfer Learning – Training on Large Dataset



Step 2: Take your custom data and **initialize** the network with weights trained in Step 1



Replace last layer with new fully-connected for output nodes per new category



Initializing with Pre-Trained Network



Step 3: (Continue to) train on new dataset

- Finetune: Update all parameters
- Freeze feature layer: Update only last layer weights (used when not enough data)



Replace last layer with new fully-connected for output nodes per new category



Finetuning on New Dataset



This works extremely well! It was surprising upon discovery.

- Features learned for 1000 object categories will work well for 1001st!
- Generalizes even across tasks (classification to object detection)



From: Razavian et al., CNN Features off-the-shelf: an Astounding Baseline for Recognition

Surprising Effectiveness of Transfer Learning



Learning with Less Labels

But it doesn't always work that well!

- If the source dataset you train on is very different from the target dataset, transfer learning is not as effective
- If you have enough data for the target domain, it just results in faster convergence

See He et al., "Rethinking ImageNet Pre-training"



Effectiveness of More Data



From: Revisiting the Unreasonable Effectiveness of Data https://ai.googleblog.com/2017/07/revisitingunreasonable-effectiveness.html



From: Hestness et al., Deep Learning Scaling Is *Predictable*



There is a large number of different low-labeled settings in DL research

Setting	Source	Target	Shift Type
Semi-supervised	Single labeled	Single unlabeled	None
Domain Adaptation	Single labeled	Single unlabeled	Non-semantic
Domain Generalization	Multiple labeled	Unknown	Non-semantic
Cross-Task Transfer	Single labeled	Single unlabeled	Semantic
Few-Shot Learning	Single labeled	Single few-labeled	Semantic
Un/Self-Supervised	Single unlabeled	Many labeled	Both/Task



Semantic Shift

Dealing with Low-Labeled Situations

