# Lecture 5: Convolutional Neural Networks

Ranjay Krishna



## Administrative: EdStem

Please make sure to check and read all pinned EdStem posts.





## Administrative: Assignment 1

Due 4/16 11:59pm

- K-Nearest Neighbor
- Linear classifiers: SVM, Softmax

#### Ranjay Krishna



## Administrative: Assignment 2

Will be released tomorrow

Due 4/25 11:59pm

- Multi-layer Neural Networks,
- Image Features,
- Optimizers

#### Ranjay Krishna



## Administrative: Fridays

This Friday

**Convolutions & Vectorization** 





## Administrative: Course Project

Project proposal due 4/29 11:59pm

"Is X a valid project for 493G1?"

- Anything related to deep learning or computer vision
- Maximum of 3 students per team
- Make a EdStem private post or come to TA Office Hours

More info on the website

#### Ranjay Krishna

#### Lecture 5 - 6

Last time: Neural Networks





#### Ranjay Krishna

Lecture 5 - 7



Lecture 5 -

8



Lecture 5 -

9



Lecture 5 -

10



Lecture 5 -

11



Lecture 5 -

12

![](_page_12_Figure_0.jpeg)

![](_page_13_Figure_1.jpeg)

Forward pass: Compute output def f(w0, x0, w1, x1, w2): s0 = w0 \* x0 s1 = w1 \* x1 s2 = s0 + s1 s3 = s2 + w2 L = sigmoid(s3)

grad_L = 1.0
$grad_s3 = grad_L * (1 - L) * L$
grad_w2 = grad_s3
grad_s2 = grad_s3
grad_s0 = grad_s2
grad_s1 = grad_s2
grad_w1 = grad_s1 * x1
grad_x1 = grad_s1 * w1
grad_w0 = grad_s0 * x0
grad_x0 = grad_s0 * w0

#### Ranjay Krishna

#### Lecture 5 - 14

![](_page_14_Figure_1.jpeg)

def	f(w0,	x0, w1,	x1,	w2):
S	0 = w0	* x0		
s	1 = w1	* x1		
s	2 = s0	+ s1		
s	3 = s2	+ w2		
L	= sigr	noid(s3)		

Base case grad\_L = 1.0 grad\_s3 = grad\_L \* (1 - L) \* L grad\_w2 = grad\_s3 grad\_s2 = grad\_s3 grad\_s0 = grad\_s2 grad\_s1 = grad\_s2 grad\_w1 = grad\_s1 \* x1 grad\_x1 = grad\_s1 \* w1 grad\_w0 = grad\_s0 \* x0 grad\_x0 = grad\_s0 \* w0

#### Ranjay Krishna

#### Lecture 5 - 15

Forward pass: Compute output

![](_page_15_Figure_1.jpeg)

Forward pass:
Compute output

Sigmoid

e	ef	f(v	v0,	X	),	w1,	x1,
	s0	=	w0	*	X	0	
	s1	=	w1	*	X.	1	
	s2	=	s0	+	s:	1	
ļ	s3	=	s2	+	W	2	
l	L	= 9	sigr	noi	id	(s3)	

grad_L = 1.0
$grad_s3 = grad_L * (1 - L) * L$
grad_w2 = grad_s3
grad_s2 = grad_s3
grad_s0 = grad_s2
grad_s1 = grad_s2
grad_w1 = grad_s1 * x1
grad_x1 = grad_s1 * w1
grad_w0 = grad_s0 <b>*</b> x0
grad_x0 = grad_s0 * w0

#### Ranjay Krishna

#### Lecture 5 - 16

#### April 15, 2025

w2):

![](_page_16_Figure_1.jpeg)

Forward pass: Compute output

Add gate

de	ef	f(v	v0,	x	Э,	w1,	x1,
	s0	=	w0	*	x	)	
	s1	=	w1	*	x1	L	
	s2	=	s0	+	s1	L	
	s3	=	s2	+	W2	2	
	L	= 5	sigr	no:	id(	(s3)	

$grad_L = 1.0$
<u>grad_s3 = grad_L * (1 - L) * L</u>
grad_w2 = grad_s3
grad_s2 = grad_s3
grad_s0 = grad_s2
grad_s1 = grad_s2
grad_w1 = grad_s1 * x1
grad_x1 = grad_s1 * w1
grad_w0 = grad_s0 * x0
grad_x0 = grad_s0 * w0

#### Ranjay Krishna

#### Lecture 5 - 17

#### April 15, 2025

w2):

![](_page_17_Figure_1.jpeg)

	<pre>def f(w0,</pre>	x0, w1,	x1,	w2):
	s0 = w0	* X0		
Forward pass:	s1 = w1	<b>*</b> x1		
Compute output	s2 = s0	+ s1		
Compute output	s3 = s2	+ w2		
	•			

Add gate

L

grad_L = 1.0
$grad_s3 = grad_L * (1 - L) * L$
grad_w2 = grad_s3
grad_s2 = grad_s3
grad_s0 = grad_s2
grad_s1 = grad_s2
grad_w1 = grad_s1 * x1
grad_x1 = grad_s1 <b>*</b> w1
grad_w0 = grad_s0 <b>*</b> x0
arad x0 = arad s0 * w0

= sigmoid(s3)

#### Ranjay Krishna

#### Lecture 5 - 18

![](_page_18_Figure_1.jpeg)

	<pre>def f(w0, x0, w1, x1, w2)</pre>	1
Forward page:	s0 = w0 * x0	
	s1 = w1 * x1	
Compute output	s2 = s0 + s1	
	s3 = s2 + w2	
	L = sigmoid(s3)	

$grad_L = 1.0$
$grad_s3 = grad_L * (1 - L) * L$
grad_w2 = grad_s3
grad_s2 = grad_s3
grad_s0 = grad_s2
grad_s1 = grad_s2
grad_w1 = grad_s1 * x1
grad_x1 = grad_s1 <b>*</b> w1
grad_w0 = grad_s0 <b>*</b> x0
grad_x0 = grad_s0 <b>*</b> w0

#### Ranjay Krishna

#### Lecture 5 - 19

![](_page_19_Figure_1.jpeg)

Forward pass: Compute output

Multiply gate

d	ef	f()	w0,	X	0,	w1,	x1,	w2):
	sØ	) =	w0	*	хØ	)		
	s1	. =	w1	*	X	L		
	s2	2 =	s0	+	s1	L		
	s3	} =	s2	+	W2	2		
	L	=	sigr	no:	id	(s3)		

$grad_L = 1.0$
grad_s3 = grad_L * (1 - L) * L
grad_w2 = grad_s3
grad_s2 = grad_s3
grad_s0 = grad_s2
grad_s1 = grad_s2
grad_w1 = grad_s1 * x1
grad_x1 = grad_s1 * w1
grad_w0 = grad_s0 <b>*</b> x0
grad_x0 = grad_s0 <b>*</b> w0

#### Ranjay Krishna

#### Lecture 5 - 20

"Flat" Backprop: Do this for assignment 2!

### Stage your forward/backward computation!

![](_page_20_Figure_2.jpeg)

#### Ranjay Krishna

Lecture 5 - 21

"Flat" Backprop: Do this for assignment 1!

#### E.g. for two-layer neural net:

```
# receive W1,W2,b1,b2 (weights/biases), X (data)
# forward pass:
h1 = #... function of X,W1,b1
scores = #... function of h1,W2,b2
loss = #... (several lines of code to evaluate Softmax loss)
# backward pass:
dscores = #...
dh1, dW2, db2 = #...
dW1,db1 = #...
```

#### Ranjay Krishna

#### Lecture 5 - 22

### **Backprop Implementation: Modularized API**

![](_page_22_Figure_1.jpeg)

#### Graph (or Net) object (rough pseudo code)

![](_page_22_Figure_3.jpeg)

#### Ranjay Krishna

#### Lecture 5 - 23

### Modularized implementation: forward / backward API

Gate / Node / Function object: Actual PyTorch code

![](_page_23_Figure_2.jpeg)

(x,y,z are scalars)

<pre>class Multiply(torch.autograd.Function):</pre>			
@staticmethod			
<pre>def forward(ctx, x, y):</pre>	Need to stash		
ctx.save_for_backward(x, y) -	some values for		
z = x * y	use in backward		
return z			
@staticmethod			
<pre>def backward(ctx, grad_z):</pre>	_ Upstream		
x, $y = ctx.saved_tensors$	gradient		
grad_x = y * grad_z # dz/dx * dL/dz	Multiply upstream		
<pre>grad_y = x * grad_z # dz/dy * dL/dz</pre>	and local gradients		
<pre>return grad_x, grad_y</pre>			

#### Ranjay Krishna

#### Lecture 5 - 24

### Example: PyTorch operators

pytorch / pytorch 🛛 👁 Watch				\star Unstar	26,770	¥ Fork	6,340
↔ Code ① Issues 2,286 1	Pull requests 561 III Projects 4	🗉 Wiki 🔟 Ins	ights				
Tree: 517c7c9861 - pytorch / aten	/ src / THNN / generic /		Create new	v file Up	load files	Find file	History
ezyang and facebook-github-bot C	anonicalize all includes in PyTorch. (#14849)			Latest c	ommit 517	c7c9 on Dec	: 8, 2018
AbsCriterion.c	Canonicalize all includes in PyTorch. (#	14849)	4 mont			nths ago	
BCECriterion.c	Canonicalize all includes in PyTorch. (#	14849)				4 mor	nths ago
ClassNLLCriterion.c	Canonicalize all includes in PyTorch. (#	14849)				4 mor	nths ago
Col2Im.c	Canonicalize all includes in PyTorch. (#	14849)				4 mor	nths ago
ELU.c	Canonicalize all includes in PyTorch. (#	14849)				4 mor	nths ago
FeatureLPPooling.c	Canonicalize all includes in PyTorch. (#	14849)				4 mor	nths ago
GatedLinearUnit.c	Canonicalize all includes in PyTorch. (#	14849)				4 mor	nths ago
HardTanh.c	Canonicalize all includes in PyTorch. (#	14849)				4 mor	nths ago
Im2Col.c	Canonicalize all includes in PyTorch. (#	14849)				4 mor	nths ago
IndexLinear.c	Canonicalize all includes in PyTorch. (#	14849)				4 mor	nths ago
LeakyReLU.c	Canonicalize all includes in PyTorch. (#	14849)				4 mor	nths ago
LogSigmoid.c	Canonicalize all includes in PyTorch. (#	14849)				4 mor	nths ago
MSECriterion.c	Canonicalize all includes in PyTorch. (#	14849)				4 mor	nths ago
MultiLabelMarginCriterion.c	Canonicalize all includes in PyTorch. (#	14849)				4 mor	nths ago
MultiMarginCriterion.c	Canonicalize all includes in PyTorch. (#	14849)				4 mor	nths ago
RReLU.c	Canonicalize all includes in PyTorch. (#	14849)				4 mor	nths ago
Sigmoid.c	Canonicalize all includes in PyTorch. (#	14849)				4 mor	nths ago
SmoothL1Criterion.c	Canonicalize all includes in PyTorch. (#	14849)				4 mor	nths ago
SoftMarginCriterion.c	Canonicalize all includes in PyTorch. (#	14849)				4 mor	nths ago
SoftPlus.c	Canonicalize all includes in PyTorch. (#	14849)				4 mor	nths ago
SoftShrink.c	Canonicalize all includes in PyTorch. (#	14849)				4 mor	nths ago
SparseLinear.c	Canonicalize all includes in PyTorch. (#	:14849)				4 mor	nths ago
SpatialAdaptiveAveragePooling.c	Canonicalize all includes in PyTorch. (#	14849)				4 mor	nths ago
SpatialAdaptiveMaxPooling.c	Canonicalize all includes in PyTorch. (#	14849)				4 mor	nths ago
SpatialAveragePooling.c	Canonicalize all includes in PvTorch. (#	14849)				4 mor	nths ago

SpatialClassNLLCriterion.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
SpatialConvolutionMM.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
SpatialDilatedConvolution.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
SpatialDilatedMaxPooling.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
SpatialFractionalMaxPooling.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
SpatialFullDilatedConvolution.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
SpatialMaxUnpooling.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
SpatialReflectionPadding.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
SpatialReplicationPadding.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
SpatialUpSamplingBilinear.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
SpatialUpSamplingNearest.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
THNN.h	Canonicalize all includes in PyTorch. (#14849)	4 months ago
Tanh.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
TemporalReflectionPadding.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
TemporalReplicationPadding.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
TemporalRowConvolution.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
TemporalUpSamplingLinear.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
TemporalUpSamplingNearest.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
VolumetricAdaptiveAveragePoolin	Canonicalize all includes in PyTorch. (#14849)	4 months ago
VolumetricAdaptiveMaxPooling.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
VolumetricAveragePooling.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
VolumetricConvolutionMM.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
VolumetricDilatedConvolution.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
VolumetricDilatedMaxPooling.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
VolumetricFractionalMaxPooling.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
VolumetricFullDilatedConvolution.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
VolumetricMaxUnpooling.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
VolumetricReplicationPadding.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
VolumetricUpSamplingNearest.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
VolumetricUpSamplingTrilinear.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
linear_upsampling.h	Implement nn.functional.interpolate based on upsample. (#8591)	9 months ago
pooling_shape.h	Use integer math to compute output size of pooling operations (#14405)	4 months ago
i unfold.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago

#### Ranjay Krishna

#### Lecture 5 - 25

```
#ifndef TH GENERIC FILE
                                                                                          PyTorch sigmoid layer
    #define TH_GENERIC_FILE "THNN/generic/Sigmoid.c"
    #else
    void THNN_(Sigmoid_updateOutput)(
                                                                 Forward
              THNNState *state,
              THTensor *input,
              THTensor *output)
                                                           \sigma(x) =
 9
      THTensor_(sigmoid)(output, input);
    void THNN_(Sigmoid_updateGradInput)(
14
              THNNState *state,
              THTensor *gradOutput,
              THTensor *gradInput,
              THTensor *output)
18
    {
19
      THNN_CHECK_NELEMENT(output, gradOutput);
      THTensor_(resizeAs)(gradInput, output);
21
      TH_TENSOR_APPLY3(scalar_t, gradInput, scalar_t, gradOutput, scalar_t, output,
22
        scalar_t z = *output_data;
        *gradInput_data = *gradOutput_data * (1. - z) * z;
23
      );
24
25
    3
                                                                                                                                         Source
    #endif
```

#### Lecture 5 - 26

![](_page_26_Figure_0.jpeg)

#### Lecture 5 - 27

![](_page_27_Figure_0.jpeg)

#### Lecture 5 - 28

### So far: backprop with scalars

### What about vector-valued functions?

![](_page_28_Picture_2.jpeg)

![](_page_28_Picture_3.jpeg)

![](_page_28_Picture_4.jpeg)

### **Recap: Vector derivatives**

### Scalar to Scalar

 $x\in \mathbb{R}, y\in \mathbb{R}$ 

Regular derivative:

 $\frac{\partial y}{\partial x} \in \mathbb{R}$ 

If x changes by a small amount, how much will y change?

#### Ranjay Krishna

![](_page_29_Picture_8.jpeg)

## **Recap: Vector derivatives**

Scalar to Scalar

Vector to Scalar

 $x \in \mathbb{R}, y \in \mathbb{R}$ 

Regular derivative:

Derivative is Gradient:

 $x \in \mathbb{R}^N, y \in \mathbb{R}$ 

 $\frac{\partial y}{\partial x} \in \mathbb{R}$ 

$$\frac{\partial y}{\partial x} \in \mathbb{R}^N \quad \left(\frac{\partial y}{\partial x}\right)_n = \frac{\partial y}{\partial x_n}$$

If x changes by a small amount, how much will y change?

For each element of x, if it changes by a small amount then how much will y change?

#### Ranjay Krishna

![](_page_30_Picture_12.jpeg)

### Remember this example from last lecture?

![](_page_31_Figure_1.jpeg)

Lecture 5 - 32

Vector to Scalar  $\begin{bmatrix} -1.00\\ -2.00 \end{bmatrix} x \in \mathbb{R}^N, y \in \mathbb{R}$  0.73

Ranjay Krishna

#### Derivative is Gradient:

$$\frac{\partial y}{\partial x} \in \mathbb{R}^N \quad \left(\frac{\partial y}{\partial x}\right)_n = \frac{\partial y}{\partial x_n} \quad \begin{bmatrix} 0.40\\ -0.60 \end{bmatrix}$$

### **Recap: Vector derivatives**

Scalar to Scalar

 $x \in \mathbb{R}, y \in \mathbb{R}$ 

Regular derivative:

 $\frac{\partial y}{\partial x} \in \mathbb{R}$ 

If x changes by a small amount, how much will y change?

Vector to Scalar

$$x \in \mathbb{R}^N, y \in \mathbb{R}$$

Derivative is Gradient:

$$\frac{\partial y}{\partial x} \in \mathbb{R}^N \quad \left(\frac{\partial y}{\partial x}\right)_n = \frac{\partial y}{\partial x_n}$$

Vector to Vector  $x \in \mathbb{R}^N, y \in \mathbb{R}^M$ 

Derivative is Jacobian:

$$\frac{\partial y}{\partial x} \in \mathbb{R}^{N \times M} \left(\frac{\partial y}{\partial x}\right)_{n,m} = \frac{\partial y_m}{\partial x_n}$$

For each element of x, if it changes by a small amount then how much will y change? For each element of x, if it changes by a small amount then how much will each element of y change?

#### Ranjay Krishna

#### Lecture 5 -

### Backprop with Vectors

![](_page_33_Figure_1.jpeg)

#### Ranjay Krishna

Lecture 5 -

34

### Backprop with Vectors

![](_page_34_Figure_1.jpeg)

#### Ranjay Krishna

Lecture 5 -

35

### Backprop with Vectors

![](_page_35_Figure_1.jpeg)

#### Ranjay Krishna Lecture 5 -

36


# much does it influence L?

#### Ranjay Krishna

#### Lecture 5 -

#### 37



#### Ranjay Krishna

Lecture 5 -

38



#### Ranjay Krishna

#### Lecture 5 -

#### 39



#### Ranjay Krishna

#### Lecture 5 -

#### 40



#### Ranjay Krishna

#### Lecture 5 -

#### 41

Gradients of variables wrt loss have same dims as the original variable



#### Ranjay Krishna

Lecture 5 -

42



#### Ranjay Krishna

Lecture 5 -



Upstream gradient

April 15, 2025

#### Ranjay Krishna

Lecture 5 -

44

9



#### Ranjay Krishna

Lecture 5 -



[0010][5]

[0000][9]

Upstream gradient

April 15, 2025

Ranjay Krishna

Lecture 5 -

46

• **5** ·

[9]

◄\_\_\_\_\_



Ranjay Krishna

Lecture 5 -

Jacobian is **sparse**: off-diagonal entries always zero! Never **explicitly** form Jacobian -- instead use **implicit** multiplication



#### Ranjay Krishna

#### Lecture 5 -

Jacobian is **sparse**: off-diagonal entries always zero! Never **explicitly** form Jacobian -- instead use **implicit** multiplication



#### Ranjay Krishna

#### Lecture 5 -



Lecture 5 -

50



Lecture 5 -



Lecture 5 -

52



Lecture 5 -

#### 53



#### Ranjay Krishna

#### Lecture 5 - 54

#### April 15, 2025

y: [N×M]

x: [N×D] [21-3] [-342] w: [D×M] [321-1] [2132] [321-2]

Ranjay Krishna

Matrix Multiply  $y_{n,m} = \sum_{d} x_{n,d} w_{d,m}$ 

Jacobians: dy/dx: [(N×D)×(N×M)] dy/dw: [(D×M)×(N×M)]

For a neural net we may have N=64, D=M=4096 Each Jacobian takes ~256 GB of memory! Must work with them implicitly!

Lecture 5 - 55

[ 5 2 17 1] dL/dy: [N×M] [ 2 3 -3 9]

[-8 1 4 6]

y: [N×M]

[13 9 -2 -6]



#### Ranjay Krishna

#### Lecture 5 - 56

#### April 15, 2025

y: [N×M]

x: [N×D] — [21-3] [-342] w: [D×M] [321-1] [2132] [321-2] Matrix Multiply  $y_{n,m} = \sum_{d} x_{n,d} w_{d,m}$ 

**Q**: What parts of y are affected by one element of x?

A:  $x_{n,d}$  affects the whole row  $y_{n,\cdot}$ 

$$\frac{\partial L}{\partial x_{n,d}} = \sum_{m} \frac{\partial L}{\partial y_{n,m}} \frac{\partial y_{n,m}}{\partial x_{n,d}}$$

N×M

#### Ranjay Krishna

#### Lecture 5 - 57

-6 x: [N×D] Matrix Multiply 2 5 **1** |-3 ]  $y_{n,m} = \sum x_{n,d} w_{d,m}$ -3 4 2] dL/dy: [N×M] w: [D×M] 2 3 - 3 9 [-8 4 6 ] 3 2 1 - 1] **Q**: What parts of y **Q**: How much 2 1 3 2] are affected by one does  $x_{n,d}$ [321-2] element of x? affect  $y_{n,m}$ ? A:  $x_{n,d}$  affects the whole row  $y_{n,\cdot}$  $\frac{\partial L}{\partial x_{n,d}} = \sum_{m} \frac{\partial L}{\partial y_{n,m}} \frac{\partial y_{n,m}}{\partial x_{n,d}}$ 

#### Ranjay Krishna

#### Lecture 5 - 58

#### April 15, 2025

N×M



#### Ranjay Krishna

#### Lecture 5 - 59

#### April 15, 2025

IN×M

-61 13 9 -2 x: [N×D] Matrix Multiply 2 1/ 5 [ 2 **1** -3 ]  $y_{n,m} = \sum x_{n,d} w_{d,m}$  $[-3 \ 4 \ 2]$ dL/dy: [N×M] 2 3 - 3 9 w: [D×M]  $[-8 \ 1 \ 4 \ 6]$ [ 3 2 1 -1] **Q**: What parts of y **Q**: How much 2 1 3 2] are affected by one does  $x_{n,d}$ [321-2] element of x? affect  $y_{n,m}$ ? A:  $x_{n,d}$  affects the A:  $w_{d,m}$  $[N \times D]$   $[N \times M]$   $[M \times D]$ whole row  $y_{n,\cdot}$  $\frac{\partial L}{\partial x} = \left(\frac{\partial L}{\partial y}\right) w^T$  $\frac{\partial L}{\partial x_{n,d}} = \sum \frac{\partial L}{\partial y_{n,m}} \frac{\partial y_{n,m}}{\partial x_{n,d}} = \sum \frac{\partial L}{\partial y_{n,m}} w_{d,m}$ 

**IN×M** 

April 15, 2025

#### Ranjay Krishna

Lecture 5 - 60

#### Lecture 5 - 61

 $\overline{\partial u}$ 

#### April 15, 2025



 $[N \times D] [N \times M] [M \times D]$ 

x: [N×D]

By similar logic:

 $[D \times M]$   $[D \times N]$   $[N \times M]$ 

 $\frac{\partial L}{\partial w} = x^T \left( \right.$ 

2 1 -3 ] -3 4 2] w: [D×M] [ 3 2 1 -1] 2 1 3 2] [ 3 2 1 - 2]

Matrix Multiply
$$y_{n,m} = \sum_{d} x_{n,d} w_{d,m}$$

5 2 dL/dy: [N×M] 2 3 - 3 9 **(-8** 1 4 6 1



N×M -6 9

These formulas are easy to remember: they are the only way to make shapes match up!

# Wrapping up: Neural Networks



#### Ranjay Krishna

#### Lecture 5 - 62

# Next: Convolutional Neural Networks



#### Ranjay Krishna

#### Lecture 5 - 63

# **Recap: Fully Connected Layer**

32x32x3 image -> stretch to 3072 x 1



#### Ranjay Krishna

#### Lecture 5 - 64

# **Fully Connected Layer**

32x32x3 image -> stretch to 3072 x 1



#### Ranjay Krishna

#### Lecture 5 - 65

32x32x3 image -> preserve spatial structure



# Main idea: only look at small patches of an image

#### Ranjay Krishna

#### Lecture 5 - 66

**Convolution Layer** 

# 32x32x3 image



## 5x5x3 filter

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

#### Ranjay Krishna

#### Lecture 5 - 67

32x32x3 image

Filters always extend the full depth of the input volume

32 32

5x5x3 filter

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

#### Ranjay Krishna

#### Lecture 5 - 68



#### Ranjay Krishna

#### Lecture 5 - 69



#### Ranjay Krishna

#### Lecture 5 - 70



#### Ranjay Krishna

#### Lecture 5 - 71



#### Ranjay Krishna

Lecture 5 - 72
### **Convolution Layer**



#### Ranjay Krishna

Lecture 5 - 73

### **Convolution Layer**



#### activation map



#### Ranjay Krishna

#### Lecture 5 - 74

### **Convolution Layer**

#### consider a second, green filter



Ranjay Krishna

#### Lecture 5 - 75

#### April 15, 2025

28

28

For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:



We stack these up to get a "new image" of size 28x28x6!

#### Ranjay Krishna

#### Lecture 5 - 76

**Preview:** ConvNet is a sequence of Convolution Layers, interspersed with activation functions



#### Ranjay Krishna

#### Lecture 5 - 77

**Preview:** ConvNet is a sequence of Convolution Layers, interspersed with activation functions



#### Ranjay Krishna

Lecture 5 - 78

Preview

[Zeiler and Fergus 2013]

Visualization of VGG-16 by Lane McIntosh. VGG-16 architecture from [Simonyan and Zisserman 2014].

April 15, 2025



Ranjay Krishna



#### Ranjay Krishna

Lecture 5 - 80



#### Ranjay Krishna

#### Lecture 5 - 81

preview:



#### Ranjay Krishna

#### Lecture 5 - 82



#### Ranjay Krishna

#### Lecture 5 - 83



## 7x7 input (spatially) assume 3x3 filter

#### Ranjay Krishna

#### Lecture 5 - 84



## 7x7 input (spatially) assume 3x3 filter

#### Ranjay Krishna

#### Lecture 5 - 85



## 7x7 input (spatially) assume 3x3 filter

#### Ranjay Krishna

#### Lecture 5 - 86



## 7x7 input (spatially) assume 3x3 filter

#### Ranjay Krishna

#### Lecture 5 - 87



7x7 input (spatially) assume 3x3 filter

=> 5x5 output

#### Ranjay Krishna

Lecture 5 - 88



7x7 input (spatially) assume 3x3 filter applied **with stride 2** 

#### Ranjay Krishna

Lecture 5 - 89



7x7 input (spatially) assume 3x3 filter applied **with stride 2** 

#### Ranjay Krishna

Lecture 5 - 90



7x7 input (spatially) assume 3x3 filter applied with stride 2 => 3x3 output!

#### Ranjay Krishna

Lecture 5 - 91



7x7 input (spatially) assume 3x3 filter applied **with stride 3?** 

#### Ranjay Krishna

Lecture 5 - 92



7x7 input (spatially) assume 3x3 filter applied **with stride 3?** 

### doesn't fit! cannot apply 3x3 filter on 7x7 input with stride 3.

#### Ranjay Krishna

#### Lecture 5 - 93



Ν

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

#### Ranjay Krishna

#### Lecture 5 - 94

### In practice: Common to zero pad the border



e.g. input 7x7
3x3 filter, applied with stride 1
pad with 1 pixel border => what is the output?

(recall:) (N - F) / stride + 1

April 15, 2025

#### Ranjay Krishna

### In practice: Common to zero pad the border



e.g. input 7x7
3x3 filter, applied with stride 1
pad with 1 pixel border => what is the output?

7x7 output!

(recall:) (N + 2P - F) / stride + 1

April 15, 2025

#### Ranjay Krishna

### In practice: Common to zero pad the border



e.g. input 7x7
3x3 filter, applied with stride 1
pad with 1 pixel border => what is the output?

#### 7x7 output!

in general, common to see CONV layers with
stride 1, filters of size FxF, and zero-padding with
(F-1)/2. (will preserve size spatially)
e.g. F = 3 => zero pad with 1
F = 5 => zero pad with 2
F = 7 => zero pad with 3

#### Ranjay Krishna

#### Lecture 5 - 97

#### Remember back to...

E.g. 32x32 input convolved repeatedly with 5x5 filters shrinks volumes spatially! (32 -> 28 -> 24 ...). Shrinking too fast is not good, doesn't work well.



#### Ranjay Krishna

Lecture 5 - 98



### Let's assume output size is HxWxD. What is D?



Ranjay Krishna





# Let's assume output size is HxWxD. What is D? **10**



#### Ranjay Krishna



Examples time:

### Let's assume output size is HxWxD. What is D? 10 What is H or W?



#### Ranjay Krishna

#### Lecture 5 - 101

Examples time:

Let's assume output size is HxWxD. What is D? 10 What is H or W? (32+2\*2-5)/1+1 = 32



Ranjay Krishna

Lecture 5 - 102

Examples time:

Let's assume output size is HxWxD. What is D? 10 What is H or W? (32+2\*2-5)/1+1 = 32So the total output size is: 32x32x10

#### Ranjay Krishna

Lecture 5 - 103





Number of parameters in this layer?





Examples time:



April 15, 2025

Number of parameters in this layer? each filter has 5\*5\*3 + 1 = 76 params (+1 for bias) => 76\*10 = 760

Ranjay Krishna

### Convolution layer: summary

Let's assume input is  $W_1 \times H_1 \times C$ Conv layer needs 4 hyperparameters:

- Number of filters K
- The filter size **F**
- The stride S
- The zero padding P

This will produce an output of  $W_2 \times H_2 \times K$  where:

- 
$$W_2 = (W_1 - F + 2P)/S + 1$$

-  $H_2^2 = (H_1 - F + 2P)/S + 1$ 

Number of parameters: F<sup>2</sup>KC and K biases

#### Ranjay Krishna

Lecture 5 - 106

### Convolution layer: summary

Let's assume input is  $W_1 \times H_1 \times C$ Conv layer needs 4 hyperparameters:

- Number of filters K
- The filter size F
- The stride S
- The zero padding P

This will produce an output of  $W_2 \times H_2 \times K$  where:

- 
$$W_2 = (W_1 - F + 2P)/S + 1$$

-  $H_2 = (H_1 - F + 2P)/S + 1$ 

Number of parameters: F<sup>2</sup>CK and K biases

#### Common settings:

K = (powers of 2, e.g. 32, 64, 128, 512)

- F = 3, S = 1, P = 1

- F = 5, S = 2, P = ? (whatever fits)

April 15, 2025

- F = 1, S = 1, P = 0

#### Ranjay Krishna

(btw, 1x1 convolution layers are very useful)



#### Ranjay Krishna

#### Lecture 5 - 108
### (btw, 1x1 convolution layers are a very useful)



#### Ranjay Krishna

#### Lecture 5 - 109

# Example: CONV layer in PyTorch

Conv2d

CLASS torch.nn.Conv2d(in\_channels, out\_channels, kernel\_size, stride=1, padding=0, dilation=1, groups=1, bias=True)

[SOURCE]

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  $(N, C_{\rm in}, H, W)$  and output  $(N, C_{\rm out}, H_{\rm out}, W_{\rm out})$  can be precisely described as:

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

where  $\star$  is the valid 2D cross-correlation operator, N is a batch size, C denotes a number of channels, H is a height of input planes in pixels, and W is width in pixels.

- stride controls the stride for the cross-correlation, a single number or a tuple.
- padding controls the amount of implicit zero-paddings on both sides for padding number of points for each dimension.
- dilation controls the spacing between the kernel points; also known as the à trous algorithm. It is harder to
  describe, but this link has a nice visualization of what dilation does.
- groups controls the connections between inputs and outputs. in\_channels and out\_channels must both be divisible by groups.For example,
  - · At groups=1, all inputs are convolved to all outputs.
  - At groups=2, the operation becomes equivalent to having two conv layers side by side, each seeing half the input channels, and producing half the output channels, and both subsequently concatenated.
  - At groups= in\_channels, each input channel is convolved with its

own set of filters, of size:  $\begin{bmatrix} C_{out} \\ C_{in} \end{bmatrix}$ .

The parameters kernel\_size, stride, padding, dilation can either be:

- a single int in which case the same value is used for the height and width dimension
- a tuple of two ints in which case, the first int is used for the height dimension, and the second int for the width dimension

PyTorch is licensed under BSD 3-clause.

Conv layer needs 4 hyperparameters:

- Number of filters **K**
- The filter size F
- The stride S
- The zero padding **P**

#### Ranjay Krishna

#### Lecture 5 - 110

# Example: CONV layer in Keras

Conv layer needs 4 hyperparameters:

- Number of filters **K**
- The filter size F
- The stride S
- The zero padding P

#### Conv2D

keras.layers.Conv2D(filters, kernel\_size, strides=(1, 1), padding='valid', data\_format=None, d:

2D convolution layer (e.g. spatial convolution over images).

This layer creates a convolution kernel that is convolved with the layer input to produce a tensor of outputs. If use\_bias is True, a bias vector is created and added to the outputs. Finally, if activation is not None, it is applied to the outputs as well.

When using this layer as the first layer in a model, provide the keyword argument input\_shape (tuple of integers, does not include the batch axis), e.g. input\_shape=(128, 128, 3) for 128x128 RGB pictures in data\_format="channels\_last".

#### Arguments

- filters: Integer, the dimensionality of the output space (i.e. the number of output filters in the convolution).
- kernel\_size: An integer or tuple/list of 2 integers, specifying the height and width of the 2D convolution window. Can be a single integer to specify the same value for all spatial dimensions.
- strides: An integer or tuple/list of 2 integers, specifying the strides of the convolution along the height and width. Can be a single integer to specify the same value for all spatial dimensions.
   Specifying any stride value != 1 is incompatible with specifying any dilation\_rate value != 1.
- padding: one of "valid" or "same" (case-insensitive). Note that "same" is slightly inconsistent across backends with strides != 1, as described here
- data\_format: A string, one of "channels\_last" or "channels\_first". The ordering of the dimensions in the inputs. "channels\_last" corresponds to inputs with shape (batch, height, width, channels) while "channels\_first" corresponds to inputs with shape (batch, channels, height, width). It defaults to the image\_data\_format value found in your Keras config file at ~/.keras/keras.json. If you never set it, then it will be "channels\_last".

Keras is licensed under the MIT license.

#### Ranjay Krishna

#### Lecture 5 - 111

### The brain/neuron view of CONV Layer



Ranjay Krishna

#### Lecture 5 - 112

### The brain/neuron view of CONV Layer





It's just a neuron with local

connectivity...

the result of taking a dot product between the filter and this part of the image (i.e. 5\*5\*3 = 75-dimensional dot product)



#### Ranjay Krishna

#### Lecture 5 - 113

#### **Receptive field**





An activation map is a 28x28 sheet of neuron outputs:

- 1. Each is connected to a small region in the input
- 2. All of them share parameters

"5x5 filter" -> "5x5 receptive field for each neuron"

#### Ranjay Krishna

#### Lecture 5 - 114

### The brain/neuron view of CONV Layer





E.g. with 5 filters, CONV layer consists of neurons arranged in a 3D grid (28x28x5)

There will be 5 different neurons all looking at the same region in the input volume

#### Ranjay Krishna

#### Lecture 5 - 115



#### Ranjay Krishna

#### Lecture 5 - 116

#### FOUR layers in total: CONV/ReLU/POOL/FC



#### Ranjay Krishna

#### Lecture 5 - 117

## **Pooling layer**

- makes the representations smaller and more manageable
- operates over each activation map independently:



#### 224x224x64

#### Ranjay Krishna

Lecture 5 - 118

## MAX POOLING

#### Single depth slice



y

max pool with 2x2 filters and stride 2



April 15, 2025

#### Ranjay Krishna

Χ

Lecture 5 - 119

## Pooling layer: summary

Let's assume input is  $W_1 \times H_1 \times C$ Conv layer needs 2 hyperparameters:

- The spatial extent F
- The stride **S**

This will produce an output of  $W_2 \times H_2 \times C$  where:

- $W_2 = (W_1 F)/S + 1$
- $H_2^{-} = (H_1 F)/S + 1$

Number of parameters: 0

#### Ranjay Krishna

#### Lecture 5 - 120

## Fully Connected Layer (FC layer)

- Contains neurons that connect to the entire input volume, as in ordinary Neural Networks



#### Ranjay Krishna

#### Lecture 5 - 121

## Summary

- ConvNets stack CONV,POOL,FC layers
- Trend towards smaller filters and deeper architectures
- Trend towards getting rid of POOL/FC layers (just CONV)
- Between 2012-2016 architectures looked like [(CONV-RELU)\*N-POOL?]\*M-(FC-RELU)\*K,SOFTMAX where N is usually up to ~5, M is large, 0 <= K <= 2.</li>
  - but recent advances such as ResNet/GoogLeNet have challenged this paradigm