# Copyright Notice

These slides are distributed under the Creative Commons License.

DeepLearning.AI makes these slides available for educational purposes. You may not use or distribute these slides for commercial purposes. You may make copies of these slides and use or distribute them for educational purposes as long as you cite DeepLearning.AI as the source of the slides.
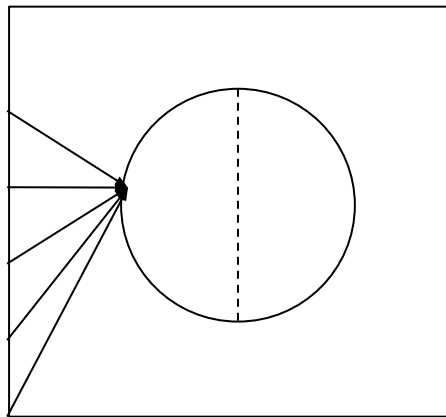
For the rest of the details of the license, see
https://creativecommons.org/licenses/by-sa/2.0/legalcode
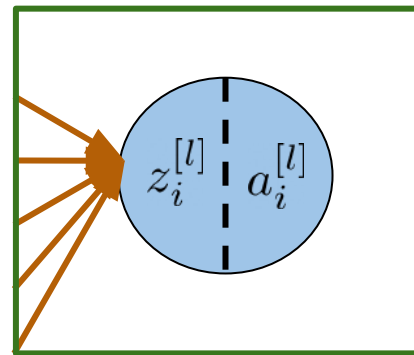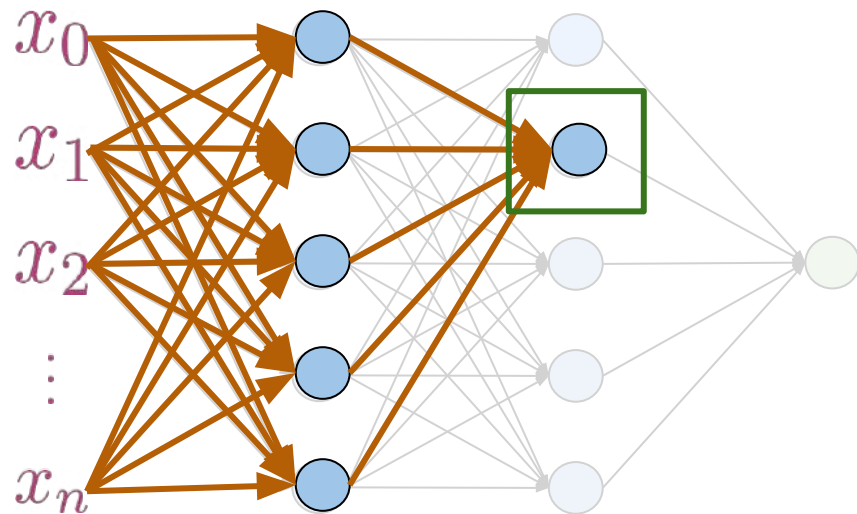
deeplearning.ai

# Activations (Basic Properties)

# Outline

- What are activations

- Reasoning behind non-linear differential activations

# Activations



$$z_i^{[l]} = \sum_{i=0} W_i^{[l]} a_i^{[l-1]}$$

$$a_i^{[l]} = \boxed{g^{[l]}}\left(z_i^{[l]}\right)$$

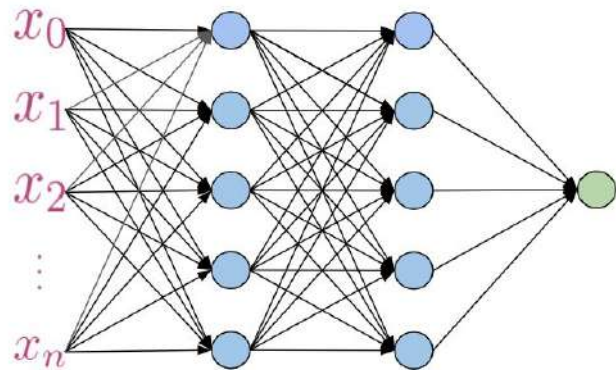Differentiable non-linear function

deeplearning.ai

# Activations

$$a_i^{[l]} = \boxed{g^{[l]}}(z_i^{[l]})$$

Differentiable non-linear function

1. Differentiable for backpropagation

2. Non-linear to compute complex features, **if not**:



$$\equiv \qquad WX + b$$

Linear regression

# Summary

- Activation functions are non-linear and differentiable

- Differentiable for backpropagation

- Non-linear to approximate complex functions

$$fx$$

deeplearning.ai

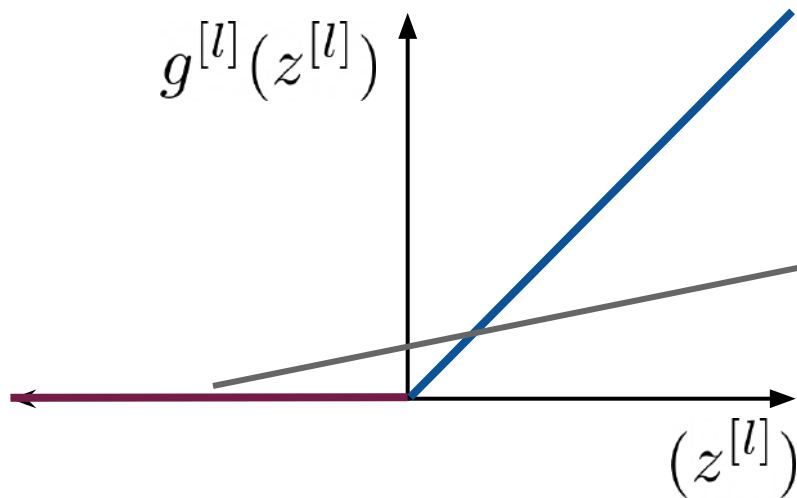Common Activation Functions

# Outline

- Common activations and their structure

  - ReLU

  - Leaky ReLU

  - Sigmoid

  - Tanh

# Activations: ReLU

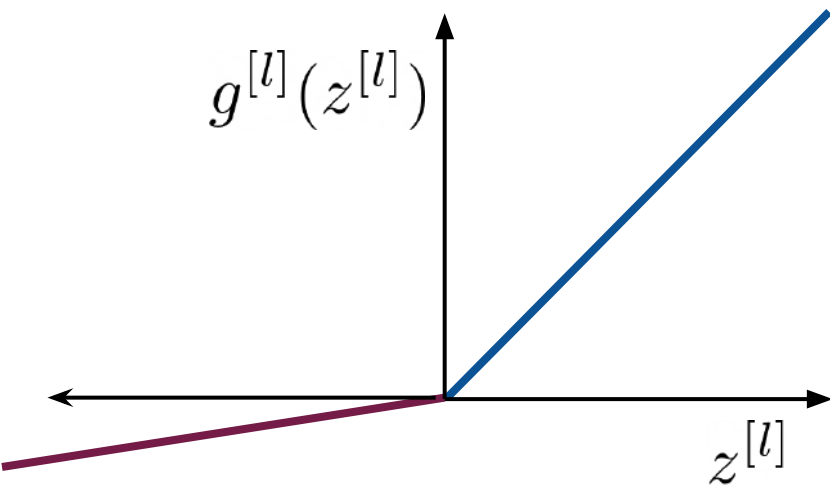ReLU = Rectified Linear Unit



$$g^{[l]}(z^{[l]}) = \max\left(\underline{0}, \underline{z^{[l]}}\right)$$

Dying ReLU problem

# Activations: Leaky ReLU
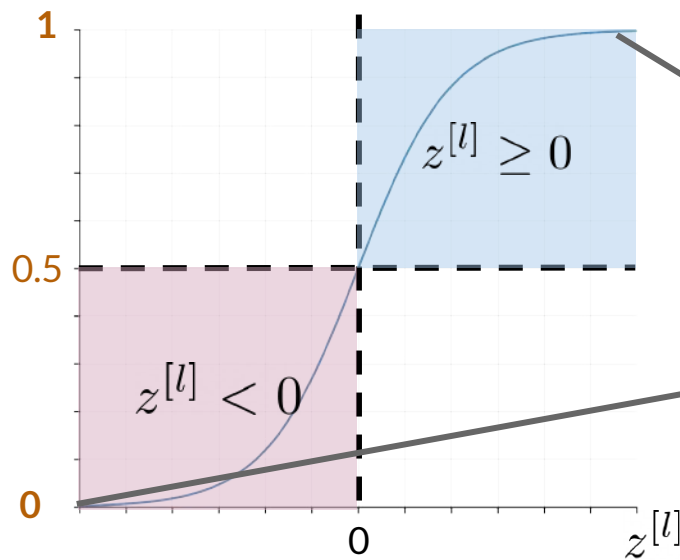


$$g^{[l]}(z^{[l]}) = max(\underline{az^{[l]}}, \underline{z^{[l]}})$$

Solves the dying
ReLU problem

# Activations: Sigmoid

$$g^{[l]}(z^{[l]}) = \frac{1}{1 + e^{-z^{[l]}}}$$



Values between 0 and 1

Vanishing gradient and saturation problems

# Activations: Tanh

$$g^{[l]}(z^{[l]}) = tanh(z^{[l]})$$



Values between -1 and 1

Keeps the sign of the input

Same issues as Sigmoid

# Summary

- ReLU activations suffer from dying ReLU

- Leaky ReLU solve the dying ReLu problem

- Sigmoid and Tanh have vanishing gradient and saturation problems

deeplearning.ai

# Batch Normalization (Explained)

# Outline

- How normalization helps models

- Internal covariate shift

- Batch normalization

# Covariate Shift

$x_1$ $w_1$

$x_2$ $w_2$

Cat

Not Cat

Change in data distribution

$x_1$

$x_2$

Cost function

$w_2$

Changes cost function

$w_1$

deeplearning.ai

# Normalization and Its Effects

# Normalization and Its Effects

$x_1$  $w_1$

$x_2$  $w_2$

Cat

Not Cat

$x_1'$   $x_2'$

**Around** mean at 0
and std. at 1

$x_1'$

$x_2'$

**Training data** uses
batch stats

**Test data** uses
training stats

# Normalization and Its Effects



$x_1 \quad w_1$

$x_2 \quad w_2$

Cat

Not Cat

$x_1' \qquad x_2'$

**Around** mean at 0 and std. at 1

Reduction of covariate shift

# Internal Covariate Shift



$z_1^{[1]}$

Changes in weights

Changes in activation distribution

# Batch Normalization



$z_1^{[1]}$

Batch Norm

Normalizes the input for each neuron

$\hat{z}_i^{[l]}$

# Summary

- Batch normalization smooths the cost function

- Batch normalization reduces the internal covariate shift

- Batch normalization speeds up learning!

# Outline

- Batch norm for training

- Batch norm for testing

# Batch Normalization: Training



$$z_i^{[l]} = \sum_{i=0}^{} W_i^{[l]} a_i^{[l-1]} \longrightarrow$$

For every example in the batch

$$\hat{z}_i^{[l]} = \frac{z_i^{[l]} - \boxed{\mu_{z_i^{[l]}}}}{\boxed{\sqrt{\sigma^2_{z_i^{[l]}} + \epsilon}}}$$

Batch mean of $z_i^{[l]}$

Batch std of $z_i^{[l]}$

deeplearning.ai

# Batch Normalization: Training



$$\hat{z}_i^{[l]} = \frac{z_i^{[l]} - \mu_{z_i^{[l]}}}{\sqrt{\sigma^2_{z_i^{[l]}} + \epsilon}}$$

$$y_i^{[l]} = \boxed{\gamma}\hat{z}_i^{[l]} + \boxed{\beta}$$

Shift factor

Scale Factor

Learnable parameters to get the optimal dist.

# Batch Normalization: Test



$$\hat{z}_i^{[l]} = \frac{z_i^{[l]} - \boxed{\mathbf{E}(z_i^{[l]})}}{\boxed{\sqrt{\mathrm{Var}(z_i^{[l]}) + \epsilon}}}$$

Running mean and running std from training

Frameworks like Tensorflow and Pytorch keep track of them

# Summary

- Batch norm introduces learnable shift and scale factors

- During test, the running statistics from training are used

- Frameworks take care of the whole process

deeplearning.ai

# Review of Convolutions

# Outline

- What convolutions are

- How they work

# What is a convolution?

# What is a convolution?



Eye Filter

# What is a convolution?



Eye Filter

Nose filter

# What is a convolution?



Eye Filter

Nose filter

Ear filter

# What is a convolution?



Eye Filter

| 5 | 3 | 1 |
|---|---|---|
| 10 | 23 | 1 |
| 1 | 42 | 4 |

Nose filter

| 2 | 21 | 5 |
|---|---|---|
| 23 | 45 | 12 |
| 3 | 32 | 4 |

Ear filter

| 32 | 2 | 24 |
|---|---|---|
| 25 | 12 | 66 |
| 3 | 45 | 2 |

# What is a convolution?

| 50 | 50 | 0 | 0 | 0 |
|----|----|---|---|---|
| 50 | 50 | 0 | 0 | 0 |
| 50 | 50 | 0 | 0 | 0 |
| 50 | 50 | 0 | 0 | 0 |
| 50 | 50 | 0 | 0 | 0 |

Grayscale image

0 (black) to 255 (white)

# What is a convolution?

| | | | | |
|---|---|---|---|---|
| 50 | 50 | 0 | 0 | 0 |
| 50 | 50 | 0 | 0 | 0 |
| 50 | 50 | 0 | 0 | 0 |
| 50 | 50 | 0 | 0 | 0 |
| 50 | 50 | 0 | 0 | 0 |

**✳**

Filter

| | | |
|---|---|---|
| 1 | 0 | -1 |
| 1 | 0 | -1 |
| 1 | 0 | -1 |

Grayscale image

0 (black) to 255 (white)

# What is a convolution?

| 50 | 50 | 0 | 0 | 0 |
|----|----|---|---|---|
| 50 | 50 | 0 | 0 | 0 |
| 50 | 50 | 0 | 0 | 0 |
| 50 | 50 | 0 | 0 | 0 |
| 50 | 50 | 0 | 0 | 0 |

Grayscale image

**Filter**

| 1 | 0 | -1 |
|---|---|----|
| 1 | 0 | -1 |
| 1 | 0 | -1 |

*

Convolution

0 (black) to 255 (white)

# What is a convolution?

| | | | | |
|---|---|---|---|---|
| 50*1 | 50*0 | 0*-1 | 0 | 0 |
| 50*1 | 50*0 | 0*-1 | 0 | 0 |
| 50*1 | 50*0 | 0*-1 | 0 | 0 |
| 50 | 50 | 0 | 0 | 0 |
| 50 | 50 | 0 | 0 | 0 |

Grayscale image

**✳**

Convolution

Filter

| 1 | 0 | -1 |
|---|---|---|
| 1 | 0 | -1 |
| 1 | 0 | -1 |

0 (black) to 255 (white)

# What is a convolution?

SUM

| | | | | |
|---|---|---|---|---|
| 50*1 | 50*0 | 0*-1 | 0 | 0 |
| 50*1 | 50*0 | 0*-1 | 0 | 0 |
| 50*1 | 50*0 | 0*-1 | 0 | 0 |
| **50** | **50** | 0 | 0 | 0 |
| **50** | **50** | 0 | 0 | 0 |

**\***

Filter

| 1 | 0 | -1 |
|---|---|---|
| 1 | 0 | -1 |
| 1 | 0 | -1 |

**=**

Result

| 150 | | |
|---|---|---|
| | | |
| | | |

Grayscale image

Convolution

0 (black) to 255 (white)

# What is a convolution?

| 50 | 50*1 | 0*0 | 0*-1 | 0 |
|----|------|-----|------|---|
| 50 | 50*1 | 0*0 | 0*-1 | 0 |
| 50 | 50*1 | 0*0 | 0*-1 | 0 |
| 50 | 50 | 0 | 0 | 0 |
| 50 | 50 | 0 | 0 | 0 |

Grayscale image

**✳**

Convolution

## Filter

| 1 | 0 | -1 |
|---|---|----|
| 1 | 0 | -1 |
| 1 | 0 | -1 |

**=**

## Result

| 150 | 150 | |
|-----|-----|---|
| | | |
| | | |

0 (black) to 255 (white)

# What is a convolution?

| | | | | |
|---|---|---|---|---|
| 50 | 50 | 0*1 | 0*0 | 0*-1 |
| 50 | 50 | 0*1 | 0*0 | 0*-1 |
| 50 | 50 | 0*1 | 0*0 | 0*-1 |
| 50 | 50 | 0 | 0 | 0 |
| 50 | 50 | 0 | 0 | 0 |

Grayscale image

$*$

Convolution

**Filter**

| 1 | 0 | -1 |
|---|---|---|
| 1 | 0 | -1 |
| 1 | 0 | -1 |

$=$

Result

| 150 | 150 | 0 |
|---|---|---|
| | | |
| | | |

0 (black) to 255 (white)

# What is a convolution?

| | | | | |
|---|---|---|---|---|
| 50 | 50 | 0 | 0 | 0 |
| 50*1 | 50*0 | 0*-1 | 0 | 0 |
| 50*1 | 50*0 | 0*-1 | 0 | 0 |
| 50*1 | 50*0 | 0*-1 | 0 | 0 |
| 50 | 50 | 0 | 0 | 0 |

Grayscale image

$*$

Convolution

**Filter**

| 1 | 0 | -1 |
|---|---|---|
| 1 | 0 | -1 |
| 1 | 0 | -1 |

$=$

Result

| 150 | 150 | 0 |
|---|---|---|
| 150 | | |
| | | |

0 (black) to 255 (white)

# What is a convolution?

| 50 | 50 | 0 | 0 | 0 |
|----|----|----|----|----|
| 50 | 50*1 | 0*0 | 0*-1 | 0 |
| 50 | 50*1 | 0*0 | 0*-1 | 0 |
| 50 | 50*1 | 0*0 | 0*-1 | 0 |
| 50 | 50 | 0 | 0 | 0 |

Grayscale image

**\***

Convolution

**Filter**

| 1 | 0 | -1 |
|---|---|----|
| 1 | 0 | -1 |
| 1 | 0 | -1 |

**=**

Result

| 150 | 150 | 0 |
|-----|-----|---|
| 150 | 150 |   |
|     |     |   |

0 (black) to 255 (white)

# What is a convolution?

| 50 | 50 | 0 | 0 | 0 |
|----|----|---|---|---|
| 50 | 50 | 0 | 0 | 0 |
| 50 | 50 | 0 | 0 | 0 |
| 50 | 50 | 0 | 0 | 0 |
| 50 | 50 | 0 | 0 | 0 |

Grayscale image

\*

Convolution

**Filter**

| 1 | 0 | -1 |
|---|---|----|
| 1 | 0 | -1 |
| 1 | 0 | -1 |

==

Result

| 150 | 150 | 0 |
|-----|-----|---|
| 150 | 150 | 0 |
| 150 | 150 | 0 |

0 (black) to 255 (white)

# Summary

- Convolutions are useful layers for processing images

- They scan the image to detect useful features

- Just element-wise products and sums!

deeplearning.ai

# Padding and Stride

# Outline

- Padding and stride

- The intuition behind padding

# Stride

| 50 | 50 | 0 | 0 | 0 |
|----|----|---|---|---|
| 50 | 50 | 0 | 0 | 0 |
| 50 | 50 | 0 | 0 | 0 |
| 50 | 50 | 0 | 0 | 0 |
| 50 | 50 | 0 | 0 | 0 |

Grayscale image

$*$

Filter

| 1 | 0 | -1 |
|---|---|----|
| 1 | 0 | -1 |
| 1 | 0 | -1 |

# Stride

| | | | | |
|---|---|---|---|---|
| 50*1 | 50*0 | 0*-1 | 0 | 0 |
| 50*1 | 50*0 | 0*-1 | 0 | 0 |
| 50*1 | 50*0 | 0*-1 | 0 | 0 |
| **50** | **50** | 0 | 0 | 0 |
| **50** | **50** | 0 | 0 | 0 |

Grayscale image

$*$

Filter

| 1 | 0 | -1 |
|---|---|---|
| 1 | 0 | -1 |
| 1 | 0 | -1 |

$=$

Result

| 150 | | |
|---|---|---|
| | | |
| | | |

# Stride

→ **1 Pixel** to the right

| 50 | 50*1 | 0*0 | 0*-1 | 0 |
|---|---|---|---|---|
| 50 | 50*1 | 0*0 | 0*-1 | 0 |
| 50 | 50*1 | 0*0 | 0*-1 | 0 |
| 50 | 50 | 0 | 0 | 0 |
| 50 | 50 | 0 | 0 | 0 |

Grayscale image

**✳**

## Filter

| 1 | 0 | -1 |
|---|---|---|
| 1 | 0 | -1 |
| 1 | 0 | -1 |

**=**

## Result

| 150 | 150 | |
|---|---|---|
| | | |
| | | |

# Stride

→ **1 Pixel** to the right

| 50 | 50 | 0*1 | 0*0 | 0*-1 |
|----|----|-----|-----|------|
| 50 | 50 | 0*1 | 0*0 | 0*-1 |
| 50 | 50 | 0*1 | 0*0 | 0*-1 |
| 50 | 50 | 0 | 0 | 0 |
| 50 | 50 | 0 | 0 | 0 |

Grayscale image

$*$

**Filter**

| 1 | 0 | -1 |
|---|---|----|
| 1 | 0 | -1 |
| 1 | 0 | -1 |

$=$

Result

| 150 | 150 | 0 |
|-----|-----|---|
|  |  |  |
|  |  |  |

# Stride

**1 Pixel** to the right

| 50 | 50 | 0 | 0 | 0 |
|----|----|---|---|---|
| 50*1 | 50*0 | 0*-1 | 0 | 0 |
| 50*1 | 50*0 | 0*-1 | 0 | 0 |
| 50*1 | 50*0 | 0*-1 | 0 | 0 |
| 50 | 50 | 0 | 0 | 0 |

**1 Pixel** down

Grayscale image

$*$

## Filter

| 1 | 0 | -1 |
|---|---|----|
| 1 | 0 | -1 |
| 1 | 0 | -1 |

$=$

## Result

| 150 | 150 | 0 |
|-----|-----|---|
| 150 |     |   |
|     |     |   |

# Stride

**1 Pixel** to the right

| | | | | |
|---|---|---|---|---|
| 50 | 50 | 0 | 0 | 0 |
| 50 | 50*1 | 0*0 | 0*-1 | 0 |
| 50 | 50*1 | 0*0 | 0*-1 | 0 |
| 50 | 50*1 | 0*0 | 0*-1 | 0 |
| 50 | 50 | 0 | 0 | 0 |

**1 Pixel** down

Grayscale image

\*

### Filter

| 1 | 0 | -1 |
|---|---|---|
| 1 | 0 | -1 |
| 1 | 0 | -1 |

=

### Result

| 150 | 150 | 0 |
|-----|-----|---|
| 150 | 150 | |
| | | |

# Stride

**1 Pixel** to the right



Grayscale image

* Filter

| 1 | 0 | -1 |
|---|---|----|
| 1 | 0 | -1 |
| 1 | 0 | -1 |

=

Result

| 150 | 150 | 0 |
|-----|-----|---|
| 150 | 150 | 0 |
| 150 | 150 | 0 |

Stride = 1

# Stride

**2 Pixels** to the right

**2 Pixels** down

| 50 | 50 | 0 | 0 | 0 |
|----|----|----|----|----|
| 50 | 50 | 0 | 0 | 0 |
| 50 | 50 | 0 | 0 | 0 |
| 50 | 50 | 0 | 0 | 0 |
| 50 | 50 | 0 | 0 | 0 |

Grayscale image

$*$

### Filter

| 1 | 0 | -1 |
|---|---|----|
| 1 | 0 | -1 |
| 1 | 0 | -1 |

$=$

### Result

Stride = 2

# Stride



**2 Pixels** to the right

2 **Pixels** down

| | | |
|---|---|---|
| 50*1 | 50*0 | 0*-1 |
| 50*1 | 50*0 | 0*-1 |
| 50*1 | 50*0 | 0*-1 |

Grayscale image

**Filter**

| 1 | 0 | -1 |
|---|---|---|
| 1 | 0 | -1 |
| 1 | 0 | -1 |

Result

| 150 | |
|---|---|
| | |

Stride = 2

deeplearning.ai

# Stride

**2 Pixels** to the right

| 50 | 50 | 0*1 | 0*0 | 0*-1 |
|----|----|-----|-----|------|
| 50 | 50 | 0*1 | 0*0 | 0*-1 |
| 50 | 50 | 0*1 | 0*0 | 0*-1 |
| 50 | 50 | 0 | 0 | 0 |
| 50 | 50 | 0 | 0 | 0 |

**2 Pixels** down

Grayscale image

\*

Filter

| 1 | 0 | -1 |
|---|---|----|
| 1 | 0 | -1 |
| 1 | 0 | -1 |

=

Result

| 150 | 0 |
|-----|---|
|     |   |

Stride = 2

# Stride

**2 Pixels** to the right

| 50 | 50 | 0 | 0 | 0 |
|----|----|----|----|----|
| 50 | 50 | 0 | 0 | 0 |
| 50*1 | 50*0 | 0*-1 | 0 | 0 |
| 50*1 | 50*0 | 0*-1 | 0 | 0 |
| 50*1 | 50*0 | 0*-1 | 0 | 0 |

**2 Pixels** down

Grayscale image

$*$

Filter

| 1 | 0 | -1 |
|---|---|----|
| 1 | 0 | -1 |
| 1 | 0 | -1 |

$=$

Result

| 150 | 0 |
|-----|---|
| 150 |   |

Stride = 2

# Stride

2 Pixels to the right



Grayscale image

Filter

Result

*

=

Stride = 2

# Stride

**2 Pixels** to the right

2 **Pixels** down

| 50 | 50 | 0 | 0 | 0 |
|----|----|---|---|---|
| 50 | 50 | 0 | 0 | 0 |
| 50 | 50 | 0 | 0 | 0 |
| 50 | 50 | 0 | 0 | 0 |
| 50 | 50 | 0 | 0 | 0 |

Grayscale image

*

Filter

| 1 | 0 | -1 |
|---|---|----|
| 1 | 0 | -1 |
| 1 | 0 | -1 |

=

Result

| 150 | 0 |
|-----|---|
| 150 | 0 |

Stride = 2

# Padding

Center gets visited four times

Image \* Filter

Stride=1

Corners get visited only once

# Padding



Image

\*

Filter

Every pixel within the image gets visited the same number of times

deeplearning.ai

# Summary

- Stride determines how the filter scans the image

- Padding is like a frame on the image

- Padding gives similar importance to the edges and the center

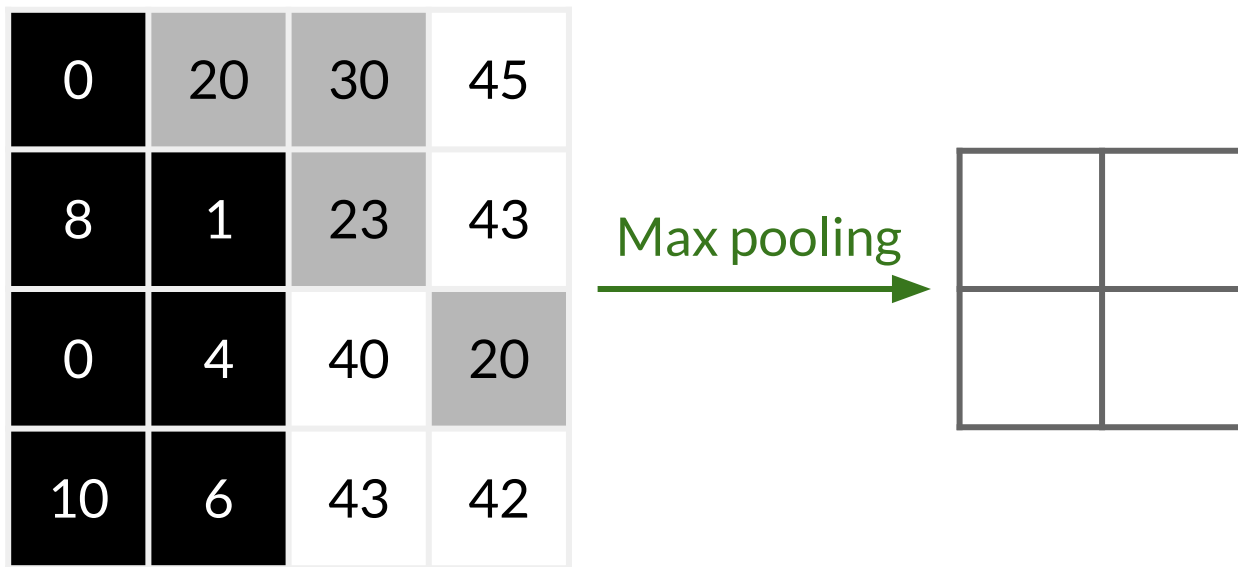deeplearning.ai

# Pooling and Upsampling

# Outline

- Pooling

- Upsampling and its relation to pooling

# Pooling



Pooling

# Max Pooling



| 0 | 20 | 30 | 45 |
|---|----|----|----|
| 8 | 1 | 23 | 43 |
| 0 | 4 | 40 | 20 |
| 10 | 6 | 43 | 42 |

Max pooling →

4x4 input to 2x2 output

deeplearning.ai

# Max Pooling



2x2 pooling with stride = 2

Max pooling →

4x4 input to 2x2 output

| 0 | 20 | 30 | 45 |
| 8 | 1 | 23 | 43 |
| 0 | 4 | 40 | 20 |
| 10 | 6 | 43 | 42 |

# Max Pooling



2x2 pooling with stride = 2

Max pooling

4x4 input to 2x2 output

# Max Pooling



2x2 pooling with stride = 2

Max pooling

4x4 input to 2x2 output

# Max Pooling



2x2 pooling with stride = 2

Max pooling

| 0 | 20 | 30 | 45 |
|---|----|----|----|
| 8 | 1 | 23 | 43 |
| 0 | 4 | 40 | 20 |
| 10 | 6 | 43 | 42 |

| 20 | 45 |
|----|----|
| 10 | |

4x4 input to 2x2 output

deeplearning.ai

# Max Pooling



2x2 pooling with stride = 2

Max pooling

4x4 input to 2x2 output

deeplearning.ai

# Max Pooling

| 0 | 20 | 30 | 45 |
|---|----|----|----|
| 8 | 1 | 23 | 43 |
| 0 | 4 | 40 | 20 |
| 10 | 6 | 43 | 42 |

Max pooling →

4x4 input to 2x2 output

2x2 pooling with stride = 2

| 20 | 45 |
|----|----|
| 10 | 43 |

Other types include:
1. Average pooling
2. Min pooling
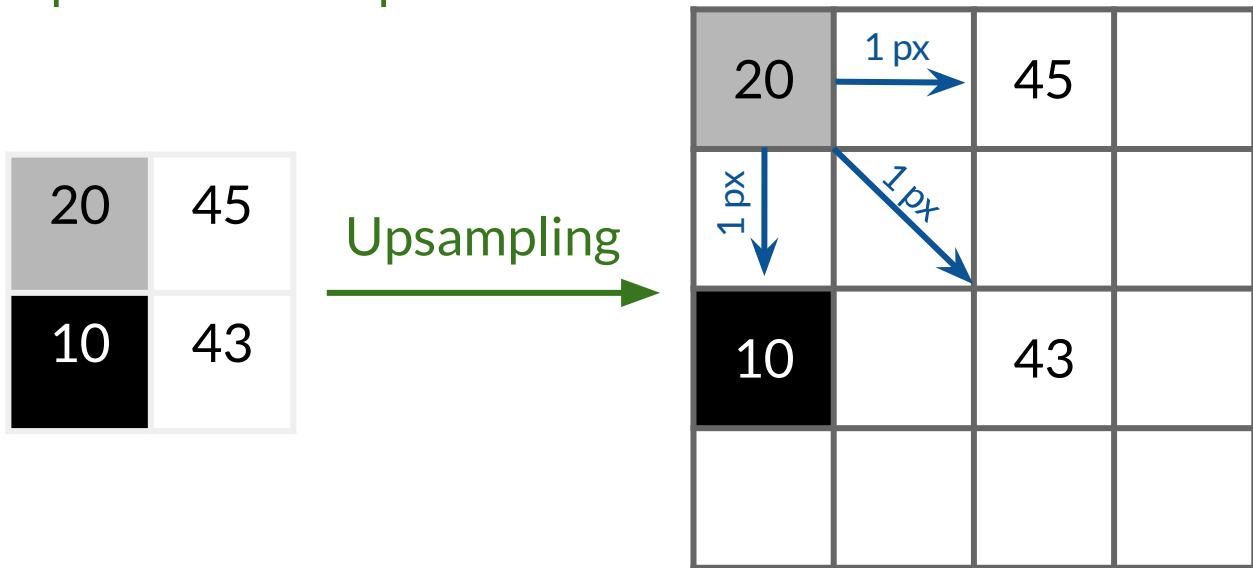
deeplearning.ai

# Upsampling



Upsampling

# Upsampling: Nearest Neighbors

2x2 input to 4x4 output

# Upsampling: Nearest Neighbors
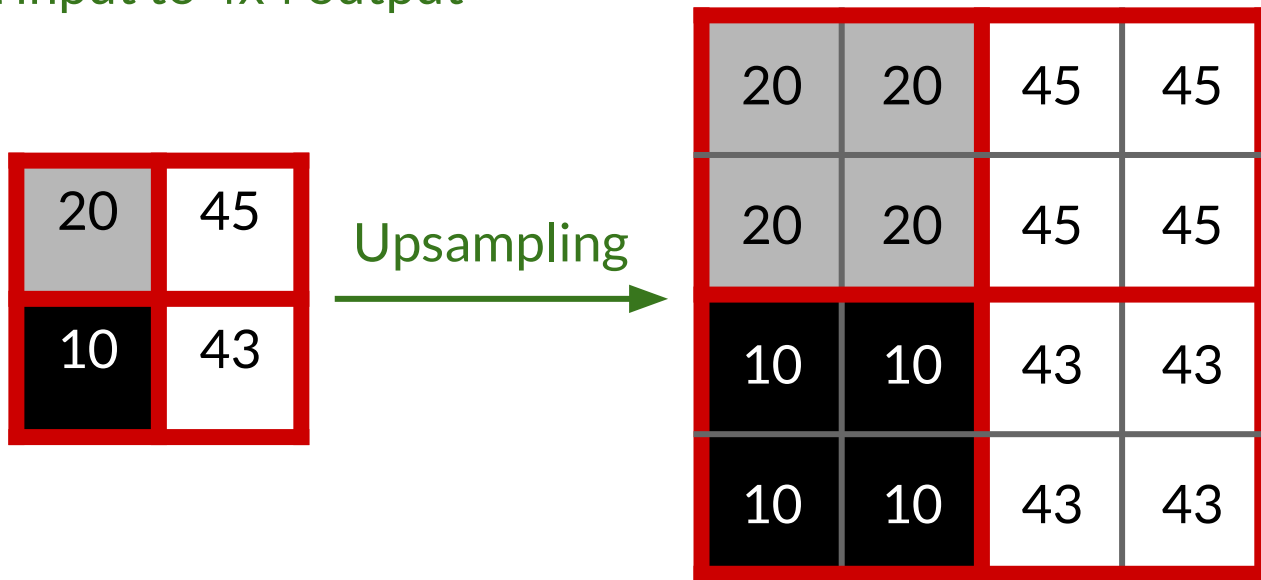
2x2 input to 4x4 output



deeplearning.ai

# Upsampling: Nearest Neighbors

2x2 input to 4x4 output

# Upsampling: Nearest Neighbors

2x2 input to 4x4 output

# Upsampling: Nearest Neighbors

2x2 input to 4x4 output



Other types include:
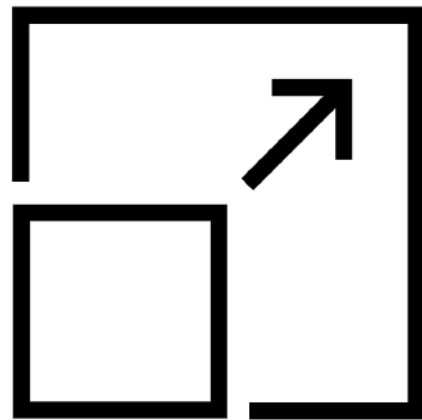1. Linear interpolation
2. Bi-linear interpolation

# Summary

- Pooling reduces the size of the input

- Upsampling increases the size of the input
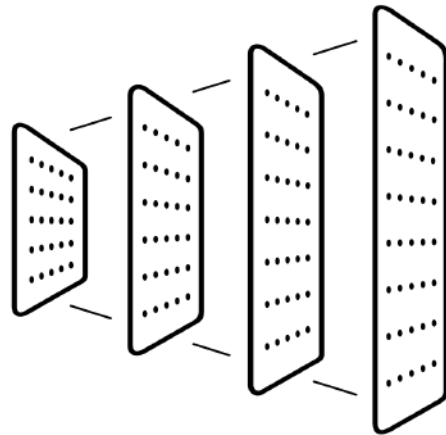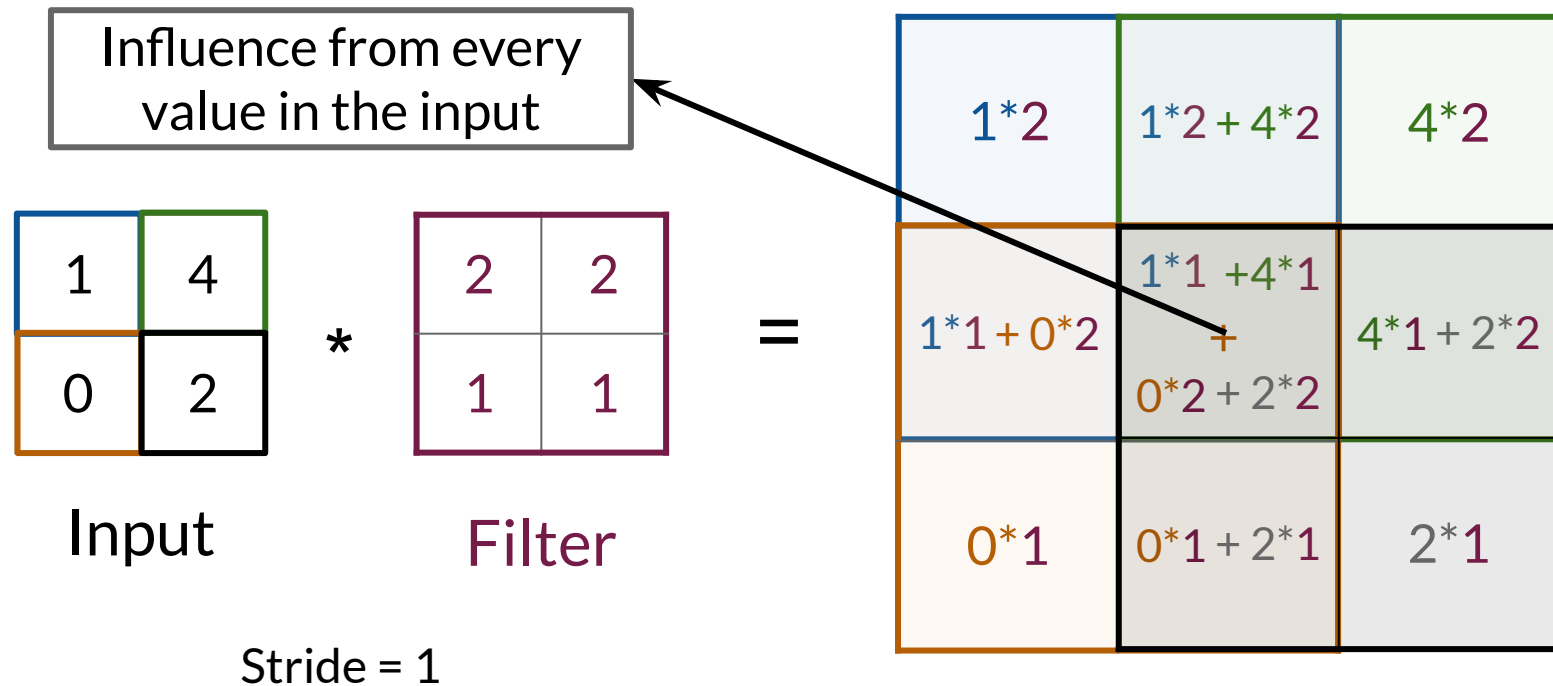
- No learnable parameters!

Transposed Convolutions

deeplearning.ai

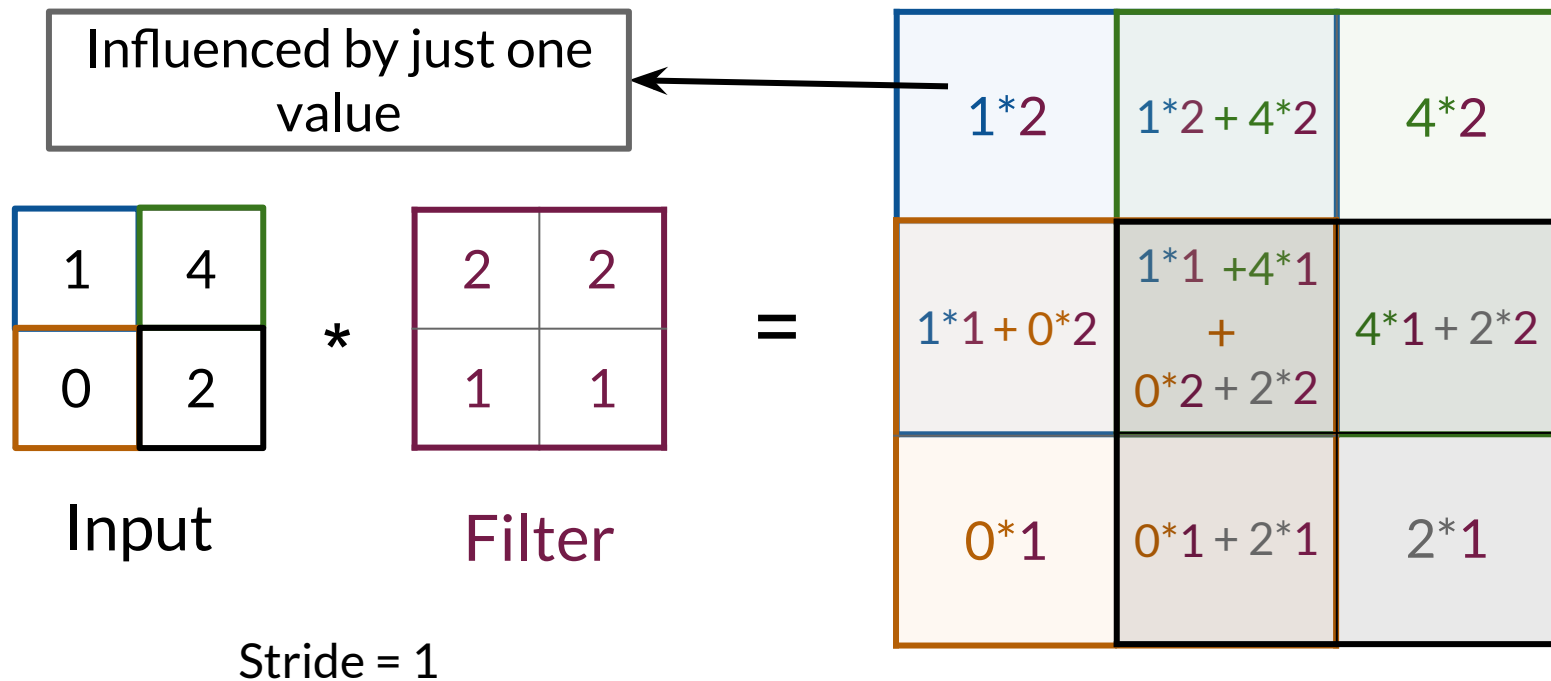# Outline

- Transposed convolutions as an upsampling technique

- Issues with transposed convolutions

# Transposed Convolution

Influence from every value in the input

| | | |
|---|---|---|
| 1*2 | 1*2 + 4*2 | 4*2 |
| 1*1 + 0*2 | 1*1 + 4*1 + 0*2 + 2*2 | 4*1 + 2*2 |
| 0*1 | 0*1 + 2*1 | 2*1 |

| 1 | 4 |
|---|---|
| 0 | 2 |

Input

\*

| 2 | 2 |
|---|---|
| 1 | 1 |

Filter

=

Stride = 1

deeplearning.ai

# Transposed Convolution

Influenced by just one value

| | | |
|---|---|---|
| 1*2 | 1*2 + 4*2 | 4*2 |
| 1*1 + 0*2 | 1*1 + 4*1 + 0*2 + 2*2 | 4*1 + 2*2 |
| 0*1 | 0*1 + 2*1 | 2*1 |

| | |
|---|---|
| 1 | 4 |
| 0 | 2 |

Input

*

| | |
|---|---|
| 2 | 2 |
| 1 | 1 |

Filter

=

Stride = 1

# The Problems with Transposed Convolution



Checkerboard Pattern

Available from: http://doi.org/10.23915/distill.00003

# Summary

- Transposed convolutions upsample

- They have learnable parameters

- Problem: results have a checkerboard pattern