# Convolutional Neural Networks

Ryan Miller

**Grinnell College**
Statistics

# Introduction

- In our previous lab we "flattened" our image data in the very first step of the neural network we created
    - This is a major limitation, as it causes us to lose the spatial information in the original organization of the pixels
- Convolutional neural networks use the mathematical operation of *convolution* to identify and extract spatially dependent hidden features
    - This is particularly effective for image data, but can be used for any data where the relative positions of input features are meaningful

**Grinnell College**
Statistics

# MNIST Example

To understand convolution, let's start with an example image from the MNIST data:



How would a standard neural network learn that this is an eight?

# Convolution

Convolution uses a *matrix of weights* known as a **kernel** (or filter) and slides it across an input matrix to generate a **feature map**:

# Stride

- **Stride** describes how a convolution kernel moves across the input to generate the feature map
  - A stride of 1 will move the kernel 1 element at a time (in the direction of rows or columns)
  - A stride of 2 will move the kernel 2 elements at a time
- A larger stride will decrease the size of the feature map, but might also result in a loss of information if an important pattern gets passed over

**Grinnell College**
Statistics

# Stride

The diagram in our earlier example uses a stride of 1. Here is the next element in the feature map:



Input tensor

Kernel

Sum up

Feature map

# Stride

The diagram below demonstrates a stride of 2. We can see that this produces a much smaller feature map:

# Multiple Input Channels

For inputs with multiple channels (such as images with RGB channels), convolution uses a separate kernel for each channel:



$$6 \times 6 \times 3 \qquad * \qquad 3 \times 3 \times 3 \qquad = \qquad 4 \times 4$$

How many values are summed to produce the (1,1) element in the 4x4 feature map?

# Remarks on Convolution

- The main benefit of convolution is that kernels are shared across multiple locations of the input
    - Thus, hidden features can be learned from different locations
- In contrast, the hidden features learned in the "vanilla" neural network architectures we've previously discussed are position-sensitive
    - For example, a neuron might learn a horizontal edge, but it could only do so for one specific location in an image
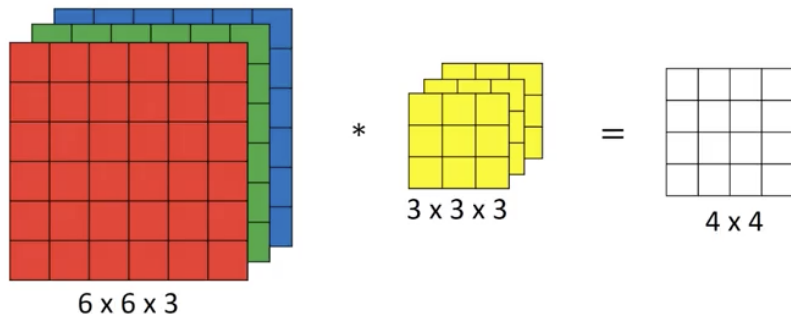    - If the data were rigorously pre-processed this might be okay, but in general it's a major limitation

**Grinnell College**
Statistics

# Example #1 - learning horizontal edges

# Example #2 - learning vertical edges

# Padding

- Convolution does not allow the center of a kernel to pass over the edges of the input tensor
  - There are some benefits to this, as it reduces the dimension of the feature map relative to the input tensor
  - But could this also cause problems?
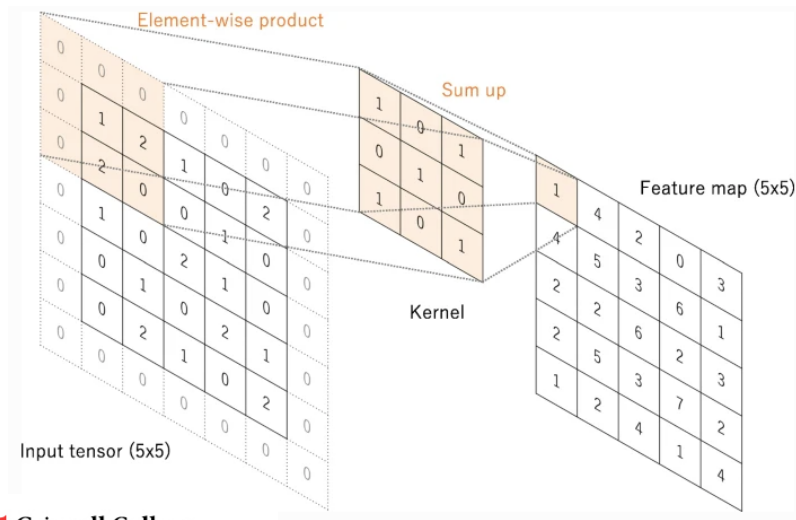
**Grinnell College**
Statistics

# Padding

- ▶ Convolution does not allow the center of a kernel to pass over the edges of the input tensor
  - ▶ There are some benefits to this, as it reduces the dimension of the feature map relative to the input tensor
  - ▶ But could this also cause problems?
- ▶ **Padding** adds extra rows and columns to create an artificial border around the input tensor
  - ▶ This can allow the center of the convolution kernel to pass over the edges of an image (or feature map)
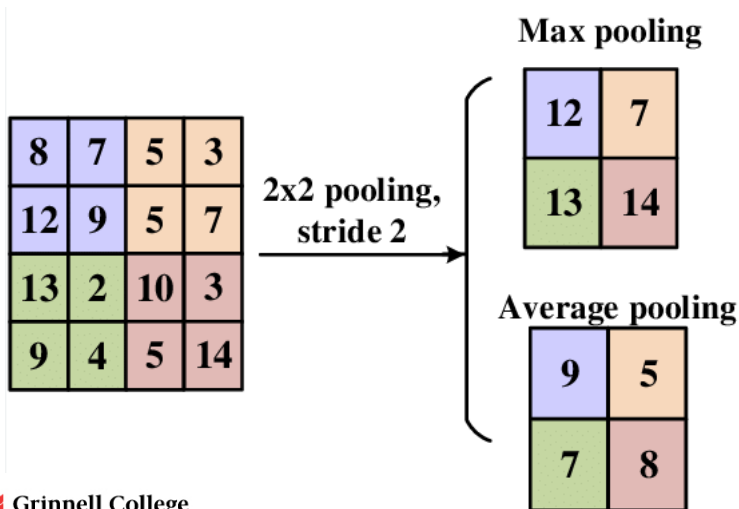  - ▶ *Zero padding*, which fills the extra rows and columns with zeros, is the most common type of padding

**Grinnell College**
Statistics

# Padding

The diagram below illustrates zero padding (of size 1):

# Pooling

- Convolutional layers can quickly increase the size and complexity of a network, especially when padding is used
- While increasing stride will reduce the size of feature maps, an operation known as **pooling** is preferable
    - *Max pooling*, which keeps the maximum value in each "patch" of an input feature map is most commonly used
    - *Average pooling*, which keeps the average value within a "patch" is sometimes used
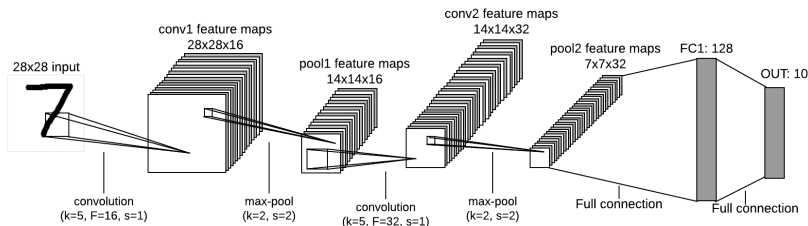
**Grinnell College**
Statistics

# Pooling

Pooling operations require a patch size and stride. Shown below is pooling using a 2x2 patch with a stride of 2:

# Architecture Diagrams

The diagram below demonstrates what the architecture of a convolutional neural network might look like:

# Remarks on Architecture

- ▶ Near the end of a convolutional neural network, the feature maps are flattened
  - ▶ It's logical to do this when the feature maps no longer contain any meaningful spatial information
  - ▶ It's also necessary to produce a properly formatted output vector
- ▶ Pooling helps keep the number of parameters in the network under control
  - ▶ For example, AlexNet contains roughly 60 million parameters despite using 3 different pooling layers to reduce complexity

**Grinnell College**
Statistics

# Image Credits

Links to images used in this presentation can be found below:

▶ "Convolutional neural networks: an overview and application in radiology"
▶ Max pooling vs. Average Pooling
▶ CNN architecture example

**Grinnell College**
Statistics