# Introduction to Convolutional Neural Networks

Ryan Miller

**Grinnell College**
Statistics

# Introduction

- ▶ The neural network architectures from our previous lecture/lab are limited by the fact that they do not account for spatial structure within the input data
- ▶ 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 is meaningful

**Grinnell College**
Statistics

# MNIST Example

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



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

# Convolution

Convolution takes 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
- Larger stride will decrease the dimensions of the feature map produced by convolution, but can also result in a loss of information (if the location where a pattern is most prominent 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:

# 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:



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

**Grinnell College**
Statistics

# 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 a specific combination pixel positions
  - If the data were rigorously pre-processed this might be okay, but in general it's a major limitation

**Grinnell College**
Statistics

# Example (learning horizontal edges)

# Example (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
  - However, if the edges of an image contain an important feature it can be problematic

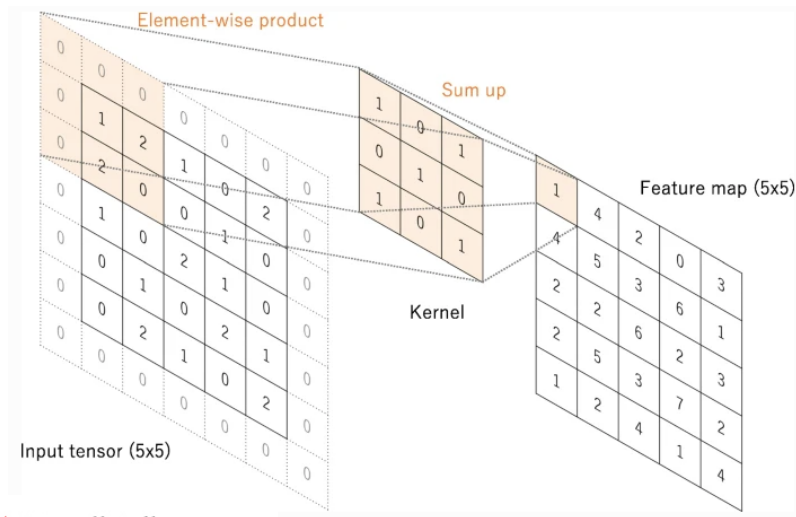**Grinnell College**
Statistics

# 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
  - However, if the edges of an image contain an important feature it can be problematic
- **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 with zeros, is the most common type of padding

**Grinnell College**
Statistics

# Padding

The diagram below illustrates zero padding:

# Pooling

- Convolutional layers can quickly increase the size and complexity of a network
- While increasing stride can help reduce the size of feature maps, an operation known as *pooling* tends to be more popular
  - *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

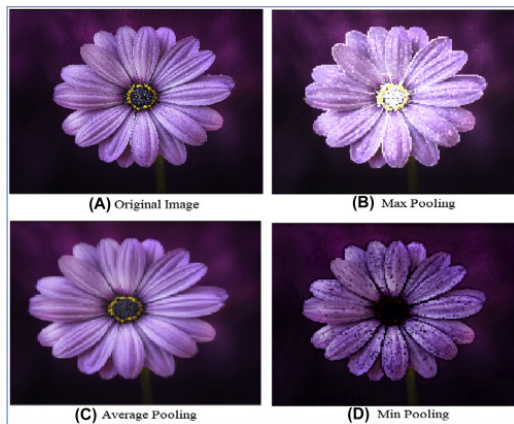**Grinnell College**
Statistics

# Pooling

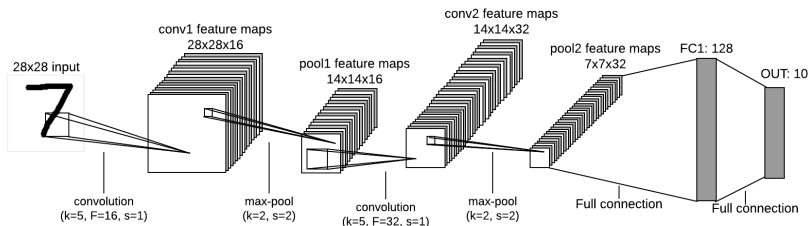Pooling operations require a patch size and stride. Shown below is pooling using a 2x2 patch with a stride of 2:

# Max vs. Average Pooling

Max pooling is widely viewed as superior because it highlights and retains the most salient spatial features:



(A) Original Image  (B) Max Pooling  (C) Average Pooling  (D) Min Pooling

# Architecture

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 map no longer contains 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 including 3 different pooling layers

**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