

Convolutional Networks: Overview

Sargur Srihari
srihari@buffalo.edu

Topics in Convolutional Networks

- Overview

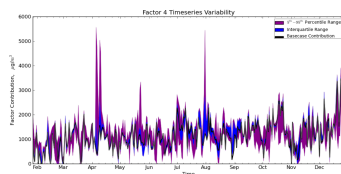
1. The Convolution Operation
2. Motivation
3. Pooling
4. Convolution and Pooling as an Infinitely Strong Prior
5. Variants of the Basic Convolution Function
6. Structured Outputs
7. Data Types
8. Efficient Convolution Algorithms
9. Random or Unsupervised Features
10. The Neuroscientific Basis for Convolutional Networks
11. Convolutional Networks and the History of Deep Learning

Plan of discussion

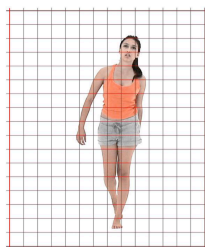
1. Overview of Convolutional Networks
2. Traditional versus Convolutional Networks
3. Topics in CNNs

Overview of Convolutional Networks

- *Convolutional networks*, also known as *Convolutional neural networks (CNNs)* are a specialized kind of neural network
- For processing data that has a known grid-like topology
 - Ex: time-series data, which is a 1-D grid, taking samples at intervals



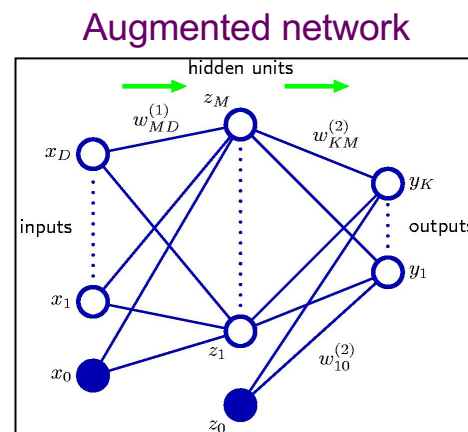
- Image data, which are 2-D grid of pixels



- They utilize convolution, which is a specialized kind of linear operation

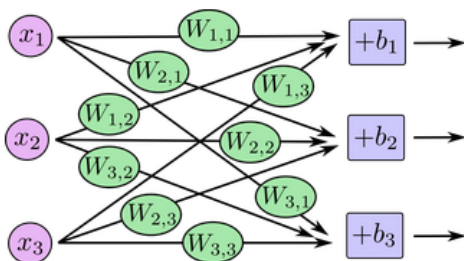
Neural net matrix multiplication

- Each layer produces values that are obtained from previous layer by performing a matrix-multiplication
 - In an unaugmented network
 - Hidden layer produces values $\mathbf{z} = h(\mathbf{W}^{(1)\top} \mathbf{x} + \mathbf{b}^{(1)})$
 - Output layer produces values $\mathbf{y} = \sigma(\mathbf{W}^{(2)\top} \mathbf{z} + \mathbf{b}^{(2)})$
 - Note: $\mathbf{W}^{(1)}$ and $\mathbf{W}^{(2)}$ are matrices rather than vectors
 - Example with $D=3, M=3$ $\mathbf{x}=[x_1, x_2, x_3]^\top$
 - We have two weight matrices $\mathbf{W}^{(1)}$ and $\mathbf{W}^{(2)}$



$$\mathbf{w} = \begin{cases} \mathbf{W}^{(1)} \text{ columns : } \mathbf{W}_1^{(1)} = [W_{11}^{(1)} W_{12}^{(1)} W_{13}^{(1)}]^\top, \mathbf{W}_2^{(1)} = [W_{21}^{(1)} W_{22}^{(1)} W_{23}^{(1)}]^\top, \mathbf{W}_3^{(1)} = [W_{31}^{(1)} W_{32}^{(1)} W_{33}^{(1)}]^\top \\ \mathbf{W}^{(2)} \text{ columns : } \mathbf{W}_1^{(2)} = [W_{11}^{(2)} W_{12}^{(2)} W_{13}^{(2)}]^\top, \mathbf{W}_2^{(2)} = [W_{21}^{(2)} W_{22}^{(2)} W_{23}^{(2)}]^\top, \mathbf{W}_3^{(2)} = [W_{31}^{(2)} W_{32}^{(2)} W_{33}^{(2)}]^\top \end{cases}$$

First Network layer



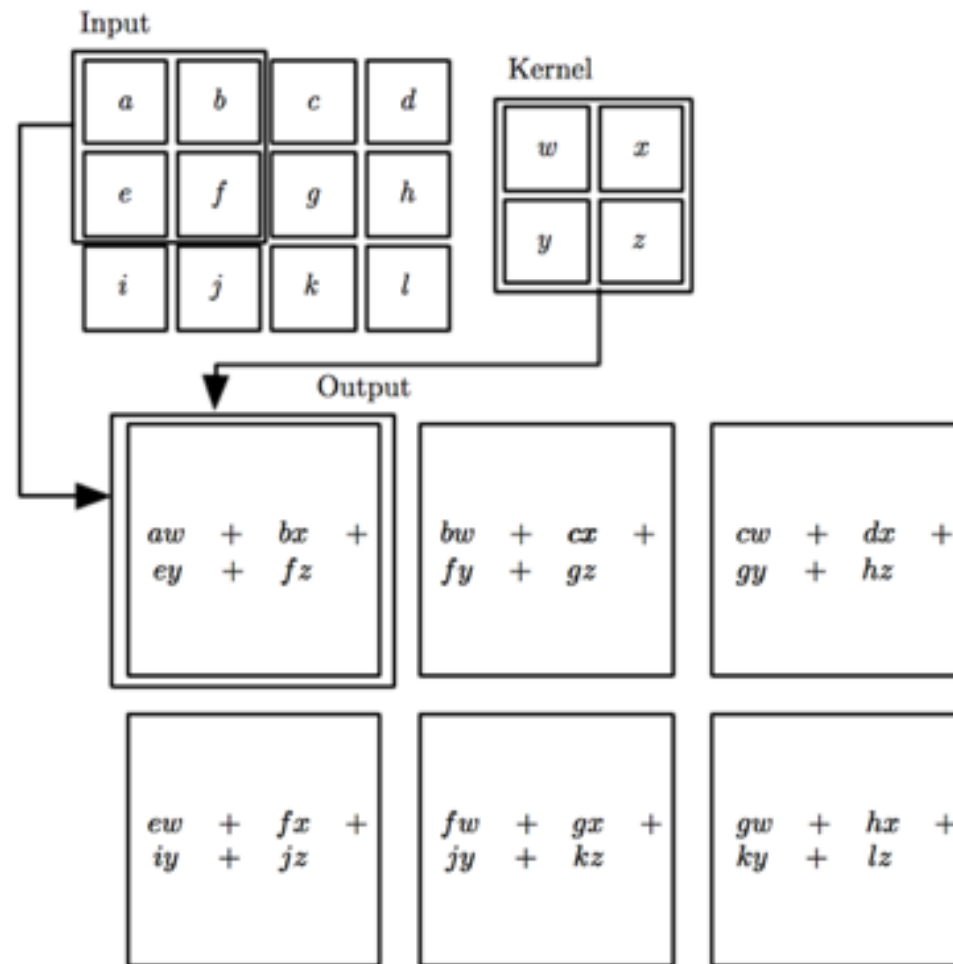
Network layer output

$$\begin{bmatrix} W_{1,1}x_1 + W_{1,2}x_2 + W_{1,3}x_3 + b_1 \\ W_{2,1}x_1 + W_{2,2}x_2 + W_{2,3}x_3 + b_2 \\ W_{3,1}x_1 + W_{3,2}x_2 + W_{3,3}x_3 + b_3 \end{bmatrix}$$

In matrix multiplication notation

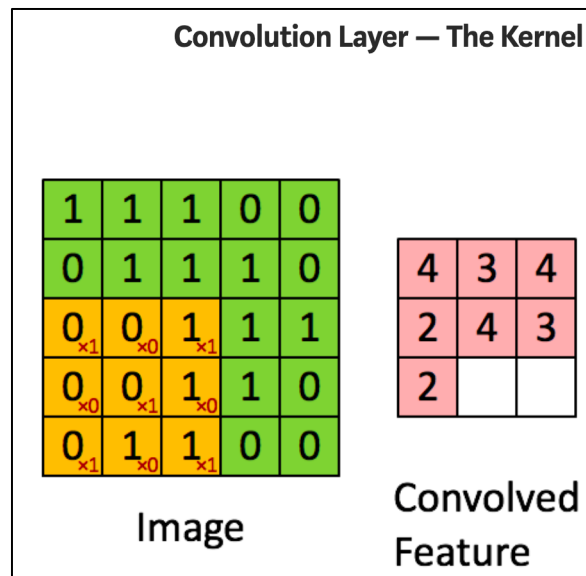
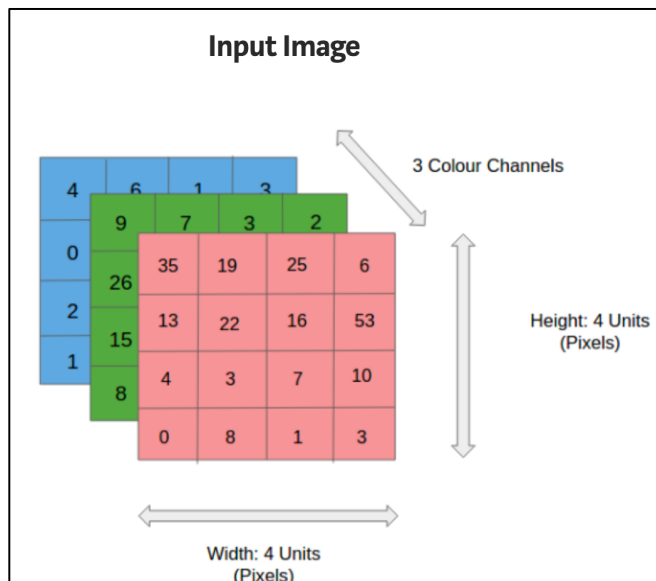
$$\begin{bmatrix} W_{1,1} & W_{1,2} & W_{1,3} \\ W_{2,1} & W_{2,2} & W_{2,3} \\ W_{3,1} & W_{3,2} & W_{3,3} \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}$$

Matrix multiplication for 2D Convolution



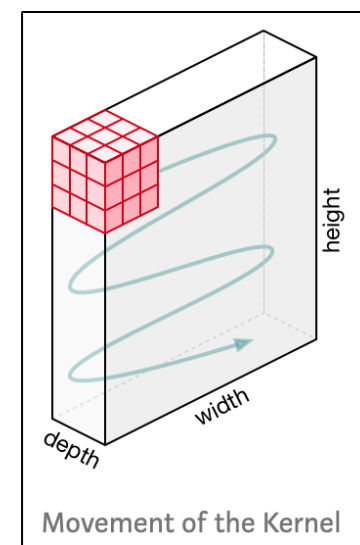
Far fewer weights needed than full matrix multiplication

Convolutional Layer for Image Recognition



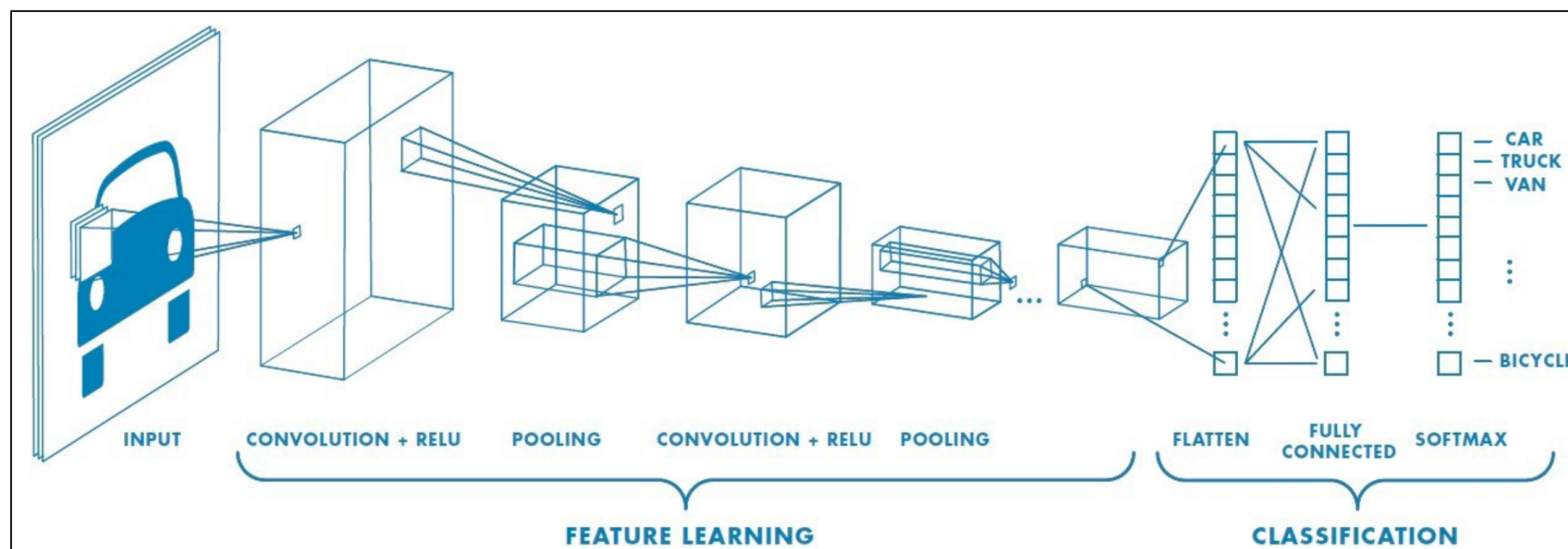
Kernel/Filter, $K =$

1	0	1
0	1	0
1	0	1



CNN is a neural network with a convolutional layer

- CNNs are simply neural networks that use convolution in place of general matrix multiplication in at least one of their layers
- Convolution can be viewed as multiplication by a matrix



Specifying a convolutional layer in Keras

- Conv1D

```
keras.layers.Conv1D(filters, kernel_size, strides=1, padding='valid', data_format='channels_last')
```

- Conv2D

```
keras.layers.Conv2D(filters, kernel_size, strides=(1, 1), padding='valid', data_format=None, dilation_rate=1)
```

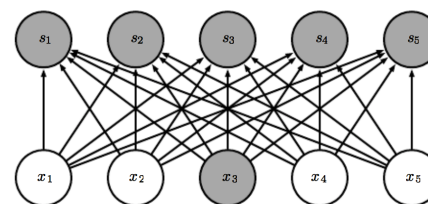
- Arguments

- **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".

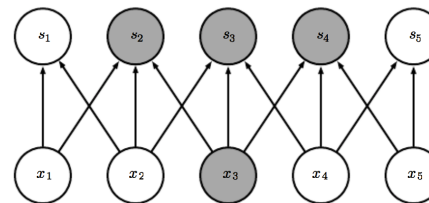
Runtime of Traditional vs Convolutional Networks

- Traditional neural network layers use matrix multiplication by a matrix of parameters with a separate parameter describing the interaction between each input unit and each output unit

$$\mathbf{s} = g(\mathbf{W}^T \mathbf{x})$$



- With m inputs and n outputs, matrix multiplication requires $m \times n$ parameters and $O(m \times n)$ runtime per example
- This means every output unit interacts with every input unit
- Convolutional network layers have sparse interactions



- If we limit no of connections for each input to k we need $k \times n$ parameters and $O(k \times n)$ runtime

Topics in Convolutional Networks

- What convolution is
- Motivation behind using convolution in a neural network
- Pooling, which almost all convolutional networks employ
- Usually the operation used in a convolutional neural network does not correspond precisely to convolution in math
 - We describe several variants on convolution function used in practice
- Making convolution more efficient
- Convolution networks stand out as an example of neuroscientific principles in deep learning
- Very deep convolutional network architectures