

Introduction to Deep Learning

Introduction to the Convolutional Network

Andres Mendez-Vazquez

June 5, 2025

Outline

1 Introduction

- The Long Path
- The Problem of Image Processing
- Multilayer Neural Network Classification
- Drawbacks
- Possible Solution

2 Convolutional Networks

- History
- Local Connectivity
- Sharing Parameters

3 Layers

- Convolutional Layer
- Convolutional Architectures
- A Little Bit of Notation

- The Reality on Models - The Use of Correlation
- Deconvolution Layer
- Non-Linearity Layer
 - Fixing the Problem, ReLu function
 - Back to the Non-Linearity Layer
- Rectification Layer
- Local Contrast Normalization Layer
- Sub-sampling and Pooling
 - Strides
- Finally, The Fully Connected Layer

4 An Example of CNN

- The Proposed Architecture
- Backpropagation
 - Deriving $w_{r,s,k}$
 - Deriving the Kernel Filters

Outline

1 Introduction

• The Long Path

- The Problem of Image Processing
- Multilayer Neural Network Classification
- Drawbacks
- Possible Solution

2 Convolutional Networks

- History
- Local Connectivity
- Sharing Parameters

3 Layers

- Convolutional Layer
- Convolutional Architectures
- A Little Bit of Notation

- The Reality on Models - The Use of Correlation
- Deconvolution Layer
- Non-Linearity Layer
- Fixing the Problem, ReLu function
- Back to the Non-Linearity Layer
- Rectification Layer
- Local Contrast Normalization Layer
- Sub-sampling and Pooling
- Strides
- Finally, The Fully Connected Layer

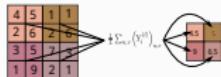
4 An Example of CNN

- The Proposed Architecture
- Backpropagation
- Deriving $w_{r,s,k}$
- Deriving the Kernel Filters

The Long Path [1]

A Small History of a Revolution

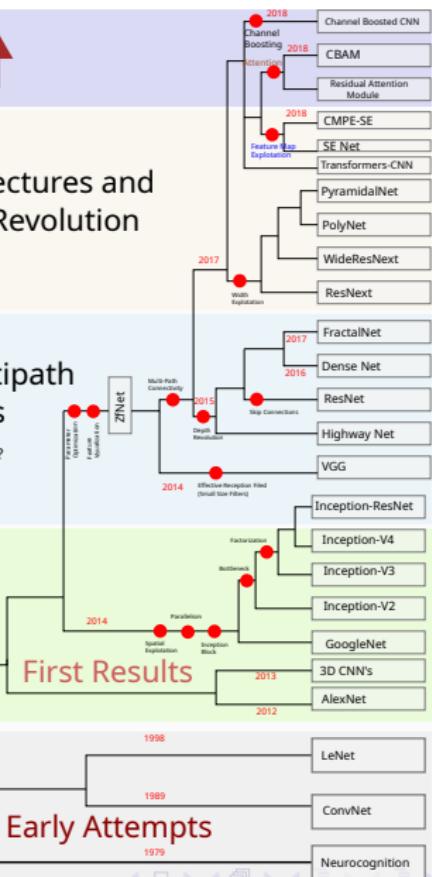
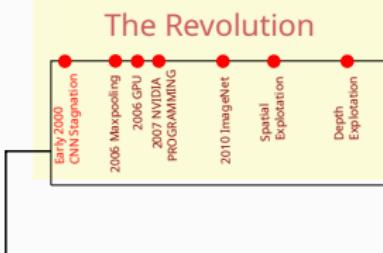
$$Y_i^{(l)}(x, y) = B_i^{(l)}(x, y) + \sum_{j=1}^{n^{(l-1)}} \frac{\sum_{k=1}^{ks}}{\sum_{s=s-k+1}^{s+k-1}} Y_j^{(l-1)}(x - s, y - l) R_{ij}^{(l)}(x, y)$$



Complex Architectures and The Attention Revolution

Residual and Multipath Architectures

The Beginng of Attention?



Outline

1 Introduction

- The Long Path
- **The Problem of Image Processing**
- Multilayer Neural Network Classification
- Drawbacks
- Possible Solution

2 Convolutional Networks

- History
- Local Connectivity
- Sharing Parameters

3 Layers

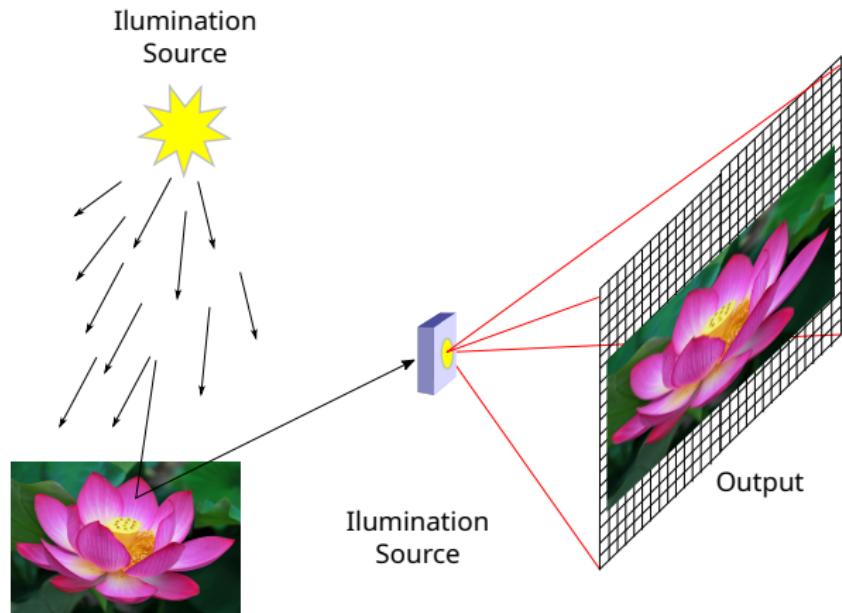
- Convolutional Layer
- Convolutional Architectures
- A Little Bit of Notation

- The Reality on Models - The Use of Correlation
- Deconvolution Layer
- Non-Linearity Layer
- Fixing the Problem, ReLu function
- Back to the Non-Linearity Layer
- Rectification Layer
- Local Contrast Normalization Layer
- Sub-sampling and Pooling
- Strides
- Finally, The Fully Connected Layer

4 An Example of CNN

- The Proposed Architecture
- Backpropagation
- Deriving $w_{r,s,k}$
- Deriving the Kernel Filters

Digital Images as pixels in a digitized matrix [2]



Further [2]

Pixel values typically represent

- Gray levels, colors, heights, opacities etc

Further [2]

Pixel values typically represent

- Gray levels, colors, heights, opacities etc

Something Notable

- Remember digitization implies that a digital image is an approximation of a real scene

Images

Common image formats include

- One sample/pixel per point (B&W or Grayscale)
- Three samples/pixel per point (Red, Green, and Blue)
- Four samples/pixel per point (Red, Green, Blue, and “Alpha”)

Therefore, we have the following process

Low Level Process

Imagen



Noise
Removal



Sharpening



Example

Edge Detection



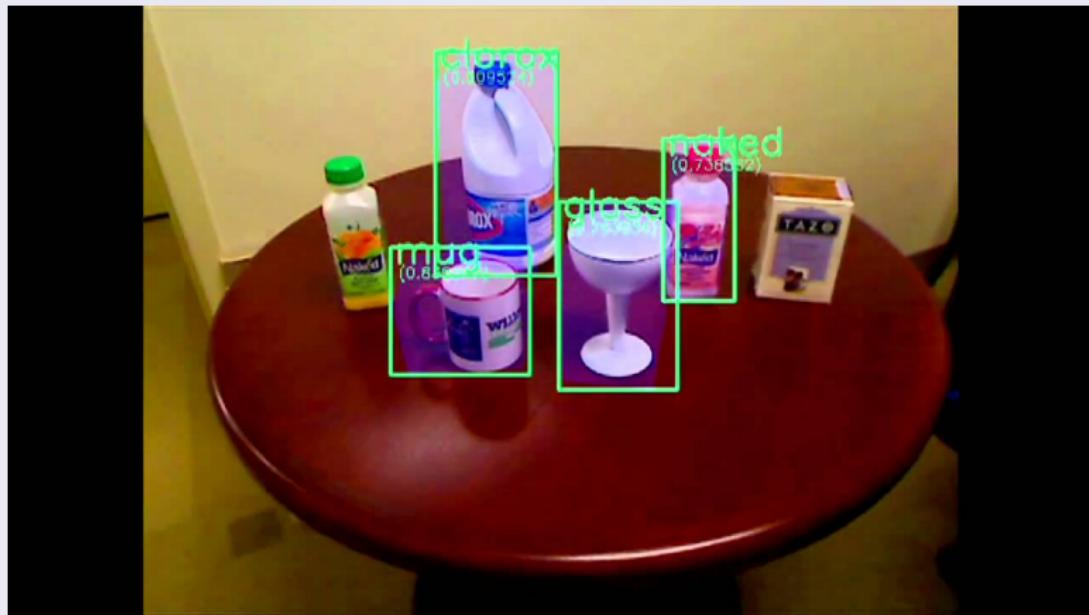
Then

Mid Level Process

Input	Processes	Output
Image	Object Recognition Segmentation	Attributes

Example

Object Recognition



Therefore

It would be nice to automatize all these processes

- We would solve a lot of headaches when setting up such process

Therefore

It would be nice to automatize all these processes

- We would solve a lot of headaches when setting up such process

Why not to use the data sets

- By using a Neural Networks that replicates the process.

Outline

1 Introduction

- The Long Path
- The Problem of Image Processing
- Multilayer Neural Network Classification
- Drawbacks
- Possible Solution

2 Convolutional Networks

- History
- Local Connectivity
- Sharing Parameters

3 Layers

- Convolutional Layer
- Convolutional Architectures
- A Little Bit of Notation

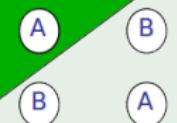
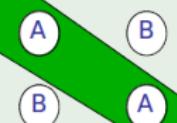
- The Reality on Models - The Use of Correlation
- Deconvolution Layer
- Non-Linearity Layer
- Fixing the Problem, ReLu function
- Back to the Non-Linearity Layer
- Rectification Layer
- Local Contrast Normalization Layer
- Sub-sampling and Pooling
- Strides
- Finally, The Fully Connected Layer

4 An Example of CNN

- The Proposed Architecture
- Backpropagation
- Deriving $w_{r,s,k}$
- Deriving the Kernel Filters

Multilayer Neural Network Classification

We have the following classification [3]

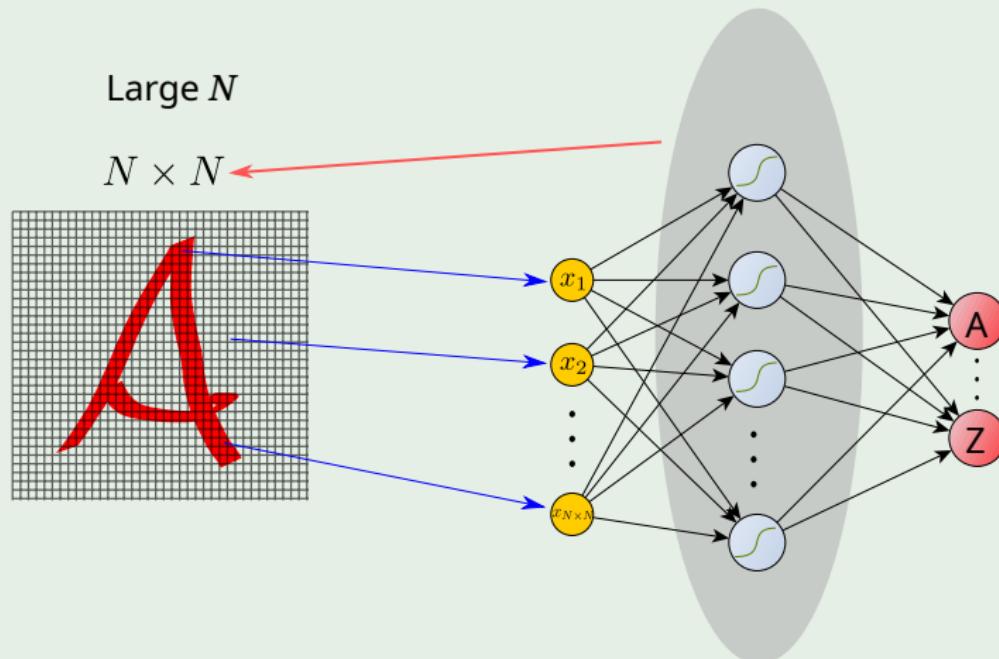
Structure	Types of Decision Regions	Exclusive-OR Problem	Classes with Meshed regions	Most General Region Shapes
Single-Layer 	Half Plane Bounded By Hyper plane			
Two-Layer 	Convex Open Or Closed Regions			
Three-Layer 	Arbitrary (Complexity Limited by No. of Nodes)			

Outline

- 1 **Introduction**
 - The Long Path
 - The Problem of Image Processing
 - Multilayer Neural Network Classification
 - **Drawbacks**
 - Possible Solution
- 2 **Convolutional Networks**
 - History
 - Local Connectivity
 - Sharing Parameters
- 3 **Layers**
 - Convolutional Layer
 - Convolutional Architectures
 - A Little Bit of Notation
- 4 **An Example of CNN**
 - The Proposed Architecture
 - Backpropagation
 - Deriving $w_{r,s,k}$
 - Deriving the Kernel Filters

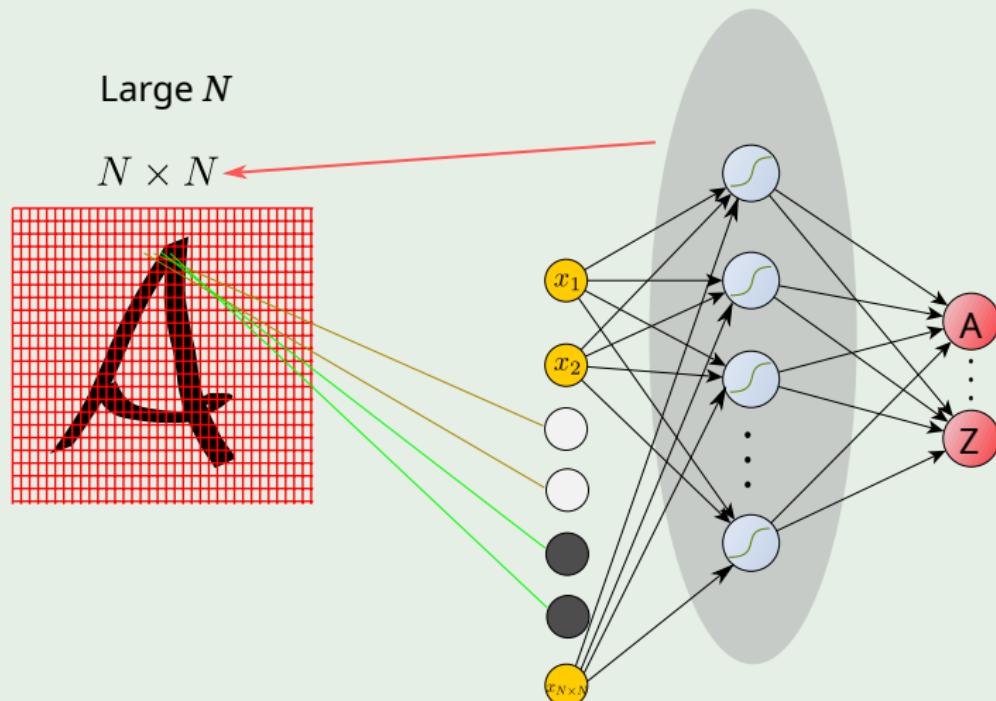
Drawbacks of previous neural networks

The number of trainable parameters becomes extremely large



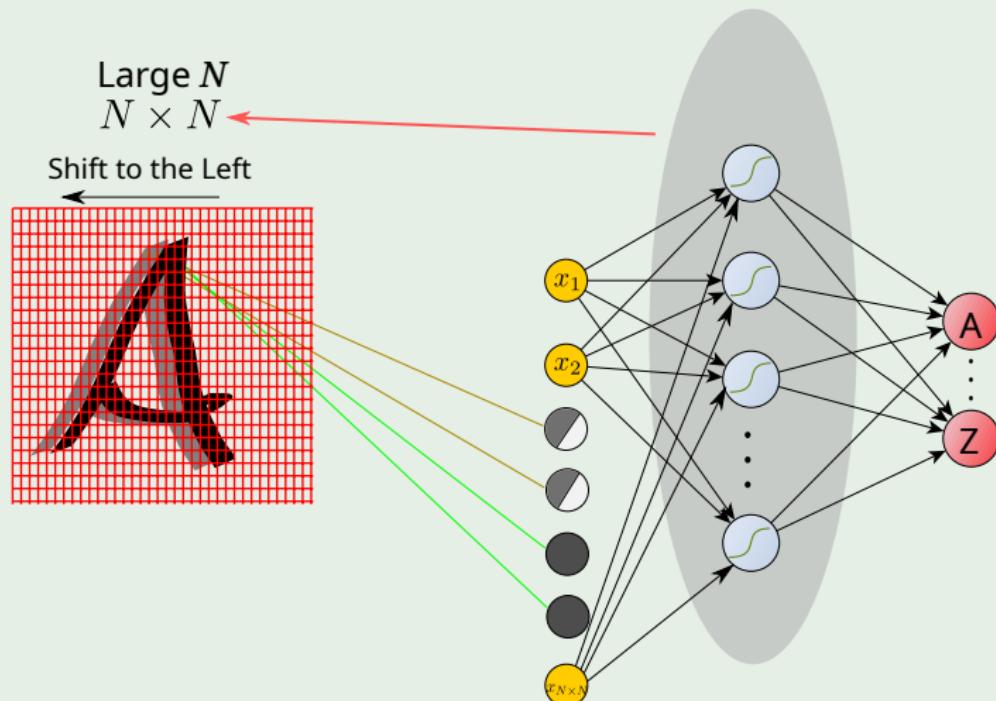
Drawbacks of previous neural networks

In addition, little or no invariance to shifting, scaling, and other forms of distortion



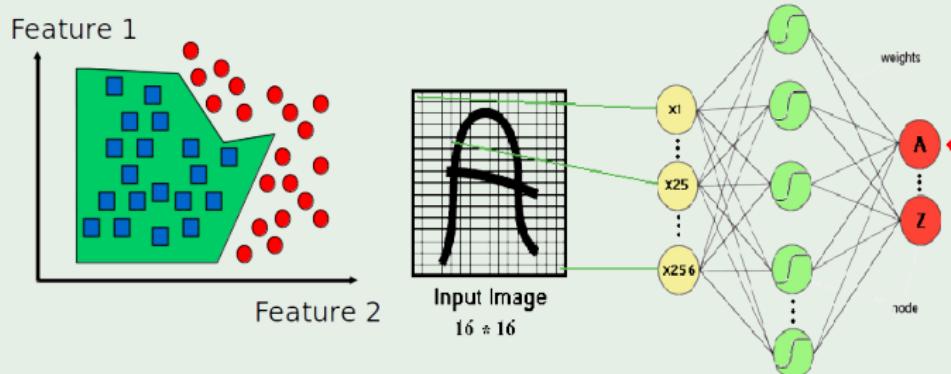
Drawbacks of previous neural networks

In addition, little or no invariance to shifting, scaling, and other forms of distortion



Drawbacks of previous neural networks

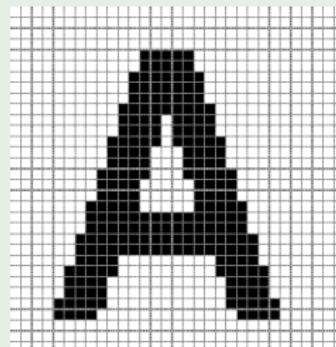
The topology of the input data is completely ignored



For Example

We have

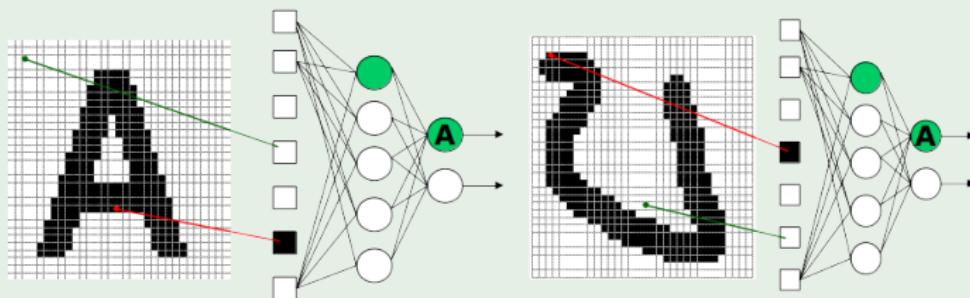
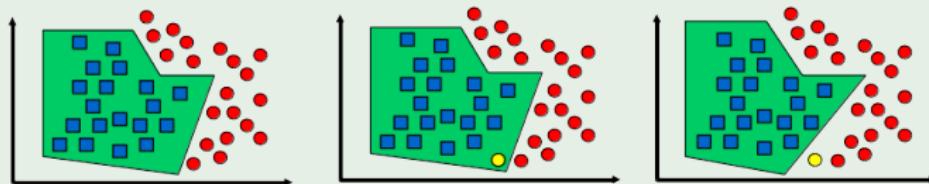
- Black and white patterns: $2^{32 \times 32} = 2^{1024}$
- Gray scale patterns: $256^{32 \times 32} = 256^{1024}$



32 * 32 input image

For Example

If we have an element that the network has never seen



Possible Solution

We can minimize this drawbacks by getting

- Fully connected network of sufficient size can produce outputs that are invariant with respect to such variations.

Possible Solution

We can minimize this drawbacks by getting

- Fully connected network of sufficient size can produce outputs that are invariant with respect to such variations.

Problem!!!

- Training time
- Network size
- Free parameters

Outline

1 Introduction

- The Long Path
- The Problem of Image Processing
- Multilayer Neural Network Classification
- Drawbacks
- Possible Solution

2 Convolutional Networks

- History
- Local Connectivity
- Sharing Parameters

3 Layers

- Convolutional Layer
- Convolutional Architectures
- A Little Bit of Notation

- The Reality on Models - The Use of Correlation
- Deconvolution Layer
- Non-Linearity Layer
- Fixing the Problem, ReLu function
- Back to the Non-Linearity Layer
- Rectification Layer
- Local Contrast Normalization Layer
- Sub-sampling and Pooling
- Strides
- Finally, The Fully Connected Layer

4 An Example of CNN

- The Proposed Architecture
- Backpropagation
- Deriving $w_{r,s,k}$
- Deriving the Kernel Filters

Hubel/Wiesel Architecture

Something Notable [4]

- D. Hubel and T. Wiesel (1959, 1962, Nobel Prize 1981)

Hubel/Wiesel Architecture

Something Notable [4]

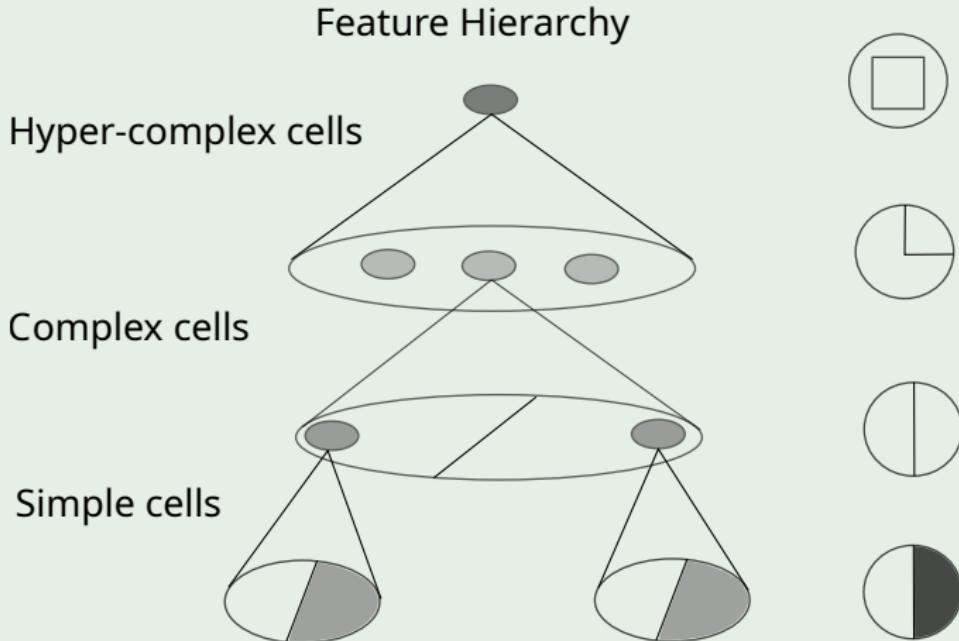
- D. Hubel and T. Wiesel (1959, 1962, Nobel Prize 1981)

They commented

- The visual cortex consists of a hierarchy of simple, complex, and hyper-complex cells

Something Like

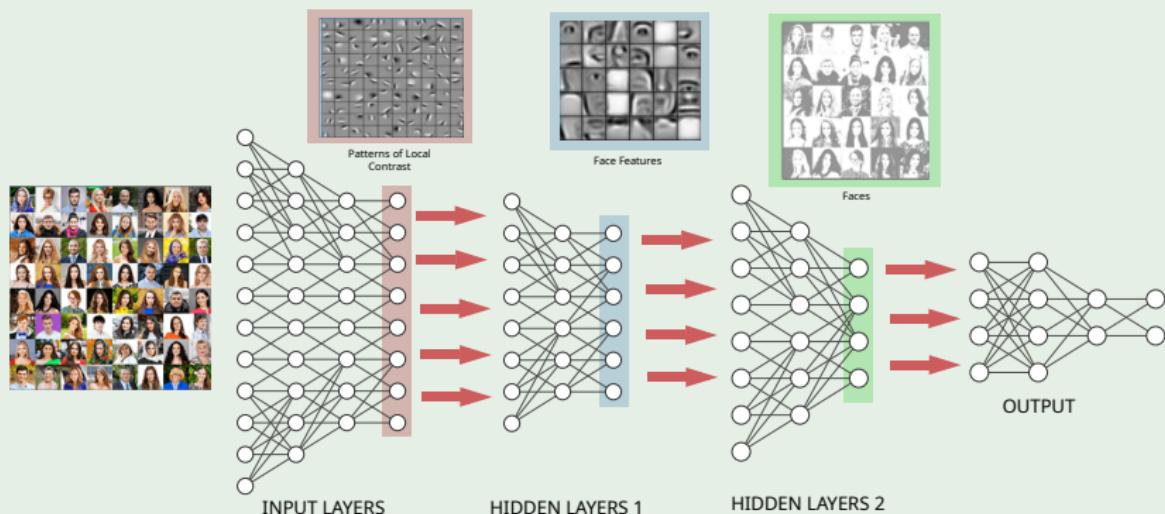
We have



History

Convolutional Neural Networks (CNN) were invented by [5]

In 1989, Yann LeCun and Yoshua Bengio introduced the concept of Convolutional Neural networks.



About CNN's

Something Notable

CNN's Were neurobiologically motivated by the findings of locally sensitive and orientation-selective nerve cells in the visual cortex.

About CNN's

Something Notable

CNN's Were neurobiologically motivated by the findings of locally sensitive and orientation-selective nerve cells in the visual cortex.

In addition

They designed a network structure that implicitly extracts relevant features.

About CNN's

Something Notable

CNN's Were neurobiologically motivated by the findings of locally sensitive and orientation-selective nerve cells in the visual cortex.

In addition

They designed a network structure that implicitly extracts relevant features.

Properties

Convolutional Neural Networks are a special kind of multi-layer neural networks.

About CNN's

In addition

- CNN is a feed-forward network that can extract topological properties from an image.

About CNN's

In addition

- CNN is a feed-forward network that can extract topological properties from an image.
- Like almost every other neural networks they are trained with a version of the back-propagation algorithm.

About CNN's

In addition

- CNN is a feed-forward network that can extract topological properties from an image.
- Like almost every other neural networks they are trained with a version of the back-propagation algorithm.
- Convolutional Neural Networks are designed to recognize visual patterns directly from pixel images with minimal preprocessing.

About CNN's

In addition

- CNN is a feed-forward network that can extract topological properties from an image.
- Like almost every other neural networks they are trained with a version of the back-propagation algorithm.
- Convolutional Neural Networks are designed to recognize visual patterns directly from pixel images with minimal preprocessing.
- They can recognize patterns with extreme variability.

Outline

1 Introduction

- The Long Path
- The Problem of Image Processing
- Multilayer Neural Network Classification
- Drawbacks
- Possible Solution

2 Convolutional Networks

- History
- Local Connectivity
- Sharing Parameters

3 Layers

- Convolutional Layer
- Convolutional Architectures
- A Little Bit of Notation

- The Reality on Models - The Use of Correlation
- Deconvolution Layer
- Non-Linearity Layer
- Fixing the Problem, ReLu function
- Back to the Non-Linearity Layer
- Rectification Layer
- Local Contrast Normalization Layer
- Sub-sampling and Pooling
- Strides
- Finally, The Fully Connected Layer

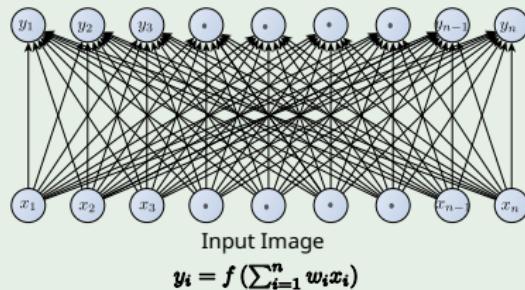
4 An Example of CNN

- The Proposed Architecture
- Backpropagation
- Deriving $w_{r,s,k}$
- Deriving the Kernel Filters

Local Connectivity

We have the following idea [6]

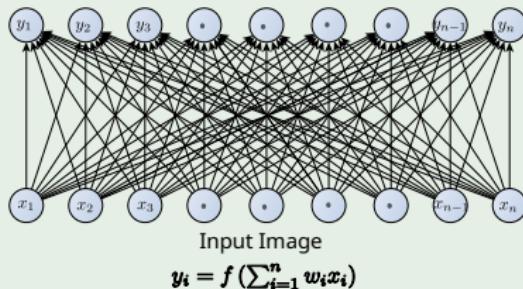
- Instead of using a full connectivity...



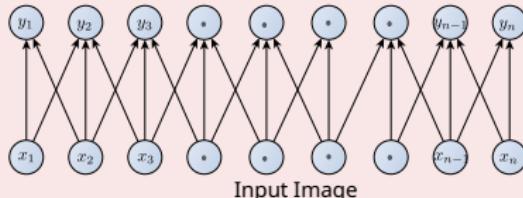
Local Connectivity

We have the following idea [6]

- Instead of using a full connectivity...



We would have something like this



Local Connectivity

We decide only to connect the neurons in a local way

- Each hidden unit is connected only to a subregion (patch) of the input image.

Local Connectivity

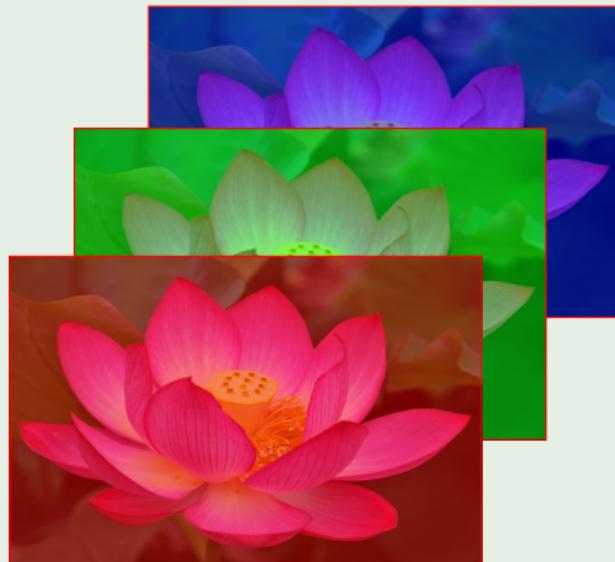
We decide only to connect the neurons in a local way

- Each hidden unit is connected only to a subregion (patch) of the input image.
- It is connected to all channels:

Local Connectivity

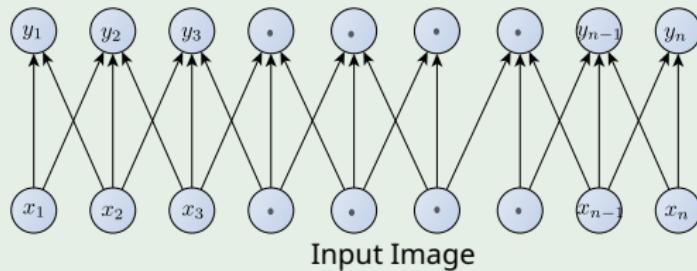
We decide only to connect the neurons in a local way

- Each hidden unit is connected only to a subregion (patch) of the input image.
- It is connected to all channels:



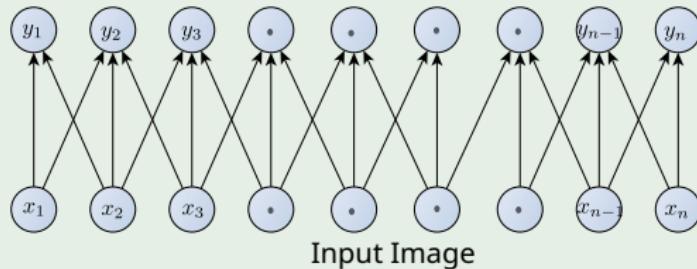
Example

For gray scale, we get something like this



Example

For gray scale, we get something like this

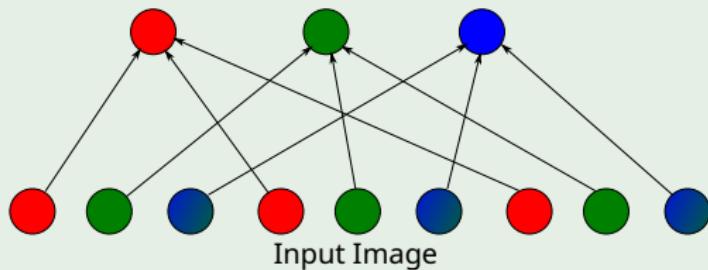


Then, our formula changes

$$y_i = f \left(\sum_{i \in L_p} w_i x_i \right) \quad (1)$$

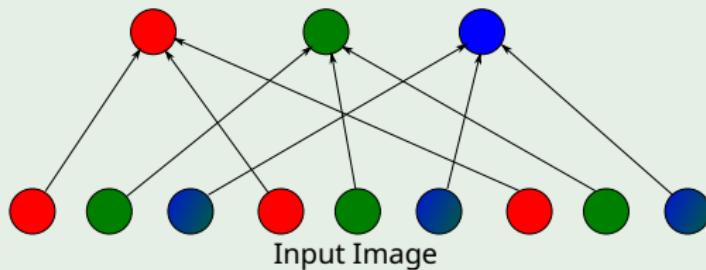
Example

In the case of the 3 channels



Example

In the case of the 3 channels



Thus

$$y_i = f \left(\sum_{i \in L_p, c} w_i x_i^c \right) \quad (2)$$

Solving the following problems...

First

- Fully connected hidden layer would have an unmanageable number of parameters

Solving the following problems...

First

- Fully connected hidden layer would have an unmanageable number of parameters

Second

- Computing the linear activation of the hidden units would have been quite expensive

How this looks in the image...

We have



Receptive Field

Outline

1 Introduction

- The Long Path
- The Problem of Image Processing
- Multilayer Neural Network Classification
- Drawbacks
- Possible Solution

2 Convolutional Networks

- History
- Local Connectivity
- Sharing Parameters

3 Layers

- Convolutional Layer
- Convolutional Architectures
- A Little Bit of Notation

- The Reality on Models - The Use of Correlation
- Deconvolution Layer
- Non-Linearity Layer
- Fixing the Problem, ReLu function
- Back to the Non-Linearity Layer
- Rectification Layer
- Local Contrast Normalization Layer
- Sub-sampling and Pooling
- Strides
- Finally, The Fully Connected Layer

4 An Example of CNN

- The Proposed Architecture
- Backpropagation
- Deriving $w_{r,s,k}$
- Deriving the Kernel Filters

Parameter Sharing

Second Idea

Share matrix of parameters across certain units.

Parameter Sharing

Second Idea

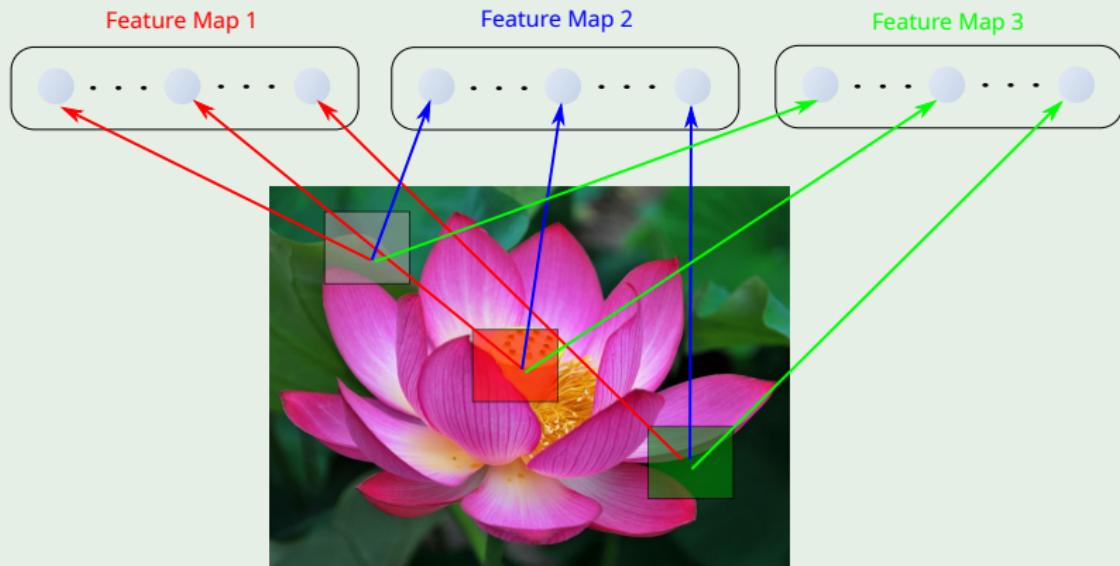
Share matrix of parameters across certain units.

These units are organized into

- The same feature “map”
 - ▶ Where the units share same parameters (For example, the same mask)

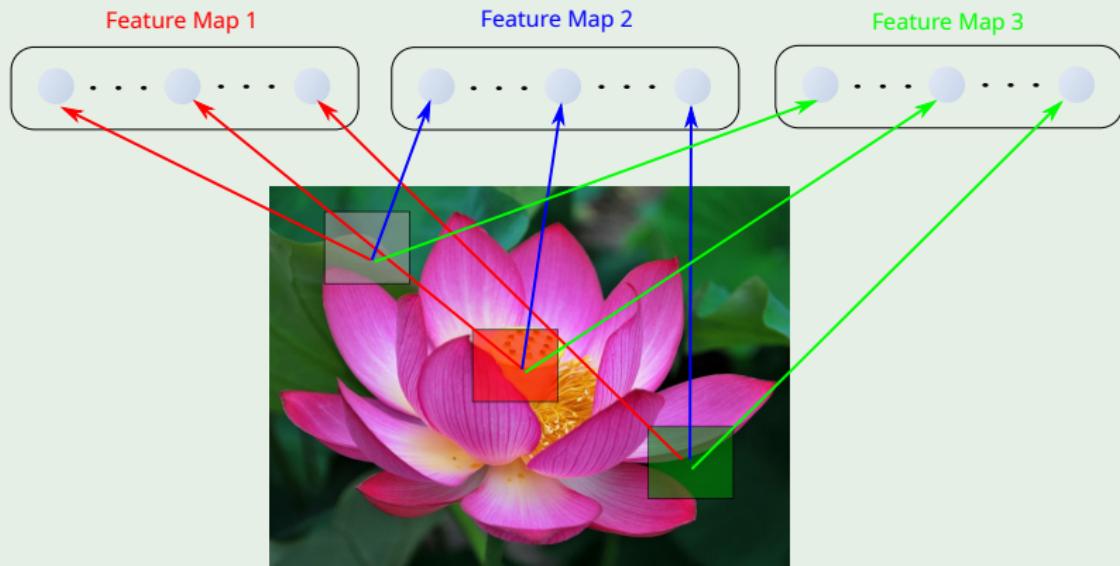
Example

We have something like this



Example

We have something like this



Now, in our notation

We have a collection of matrices representing this connectivity

- W_{ij} is the connection matrix the i th input channel with the j th feature map.

Now, in our notation

We have a collection of matrices representing this connectivity

- W_{ij} is the connection matrix the i th input channel with the j th feature map.
- In each cell of these matrices is the weight to be multiplied with the local input to the local neuron.

And now why the name convolution

Yes!!! The definition is coming now.

Outline

1 Introduction

- The Long Path
- The Problem of Image Processing
- Multilayer Neural Network Classification
- Drawbacks
- Possible Solution

2 Convolutional Networks

- History
- Local Connectivity
- Sharing Parameters

3 Layers

• Convolutional Layer

- Convolutional Architectures
- A Little Bit of Notation

- The Reality on Models - The Use of Correlation
- Deconvolution Layer
- Non-Linearity Layer
- Fixing the Problem, ReLu function
- Back to the Non-Linearity Layer
- Rectification Layer
- Local Contrast Normalization Layer
- Sub-sampling and Pooling
- Strides
- Finally, The Fully Connected Layer

4 An Example of CNN

- The Proposed Architecture
- Backpropagation
- Deriving $w_{r,s,k}$
- Deriving the Kernel Filters

Digital Images

In computer vision [2, 7]

We usually operate on digital (discrete) images:

Digital Images

In computer vision [2, 7]

We usually operate on digital (discrete) images:

- Sample the 2D space on a regular grid.

Digital Images

In computer vision [2, 7]

We usually operate on digital (discrete) images:

- Sample the 2D space on a regular grid.
- Quantize each sample (round to nearest integer).

Digital Images

In computer vision [2, 7]

We usually operate on digital (discrete) images:

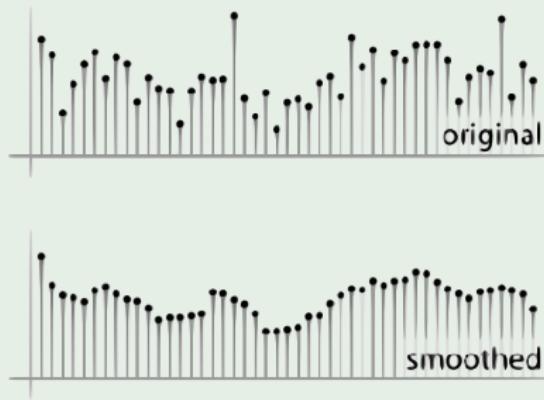
- Sample the 2D space on a regular grid.
- Quantize each sample (round to nearest integer).

The image can now be represented as a matrix of integer values,
 $I : [a, b] \times [c, d] \rightarrow [0..255]$

$$i \downarrow \begin{bmatrix} 79 & 5 & 6 & 90 & 12 & 34 & 2 & 1 \\ 8 & 90 & 12 & 34 & 26 & 78 & 34 & 5 \\ 8 & 1 & 3 & 90 & 12 & 34 & 11 & 61 \\ 77 & 90 & 12 & 34 & 200 & 2 & 9 & 45 \\ 1 & 3 & 90 & 12 & 20 & 1 & 6 & 23 \end{bmatrix} \quad j \longrightarrow$$

Many times we want to eliminate noise in a image

For example a moving average



This is defined as

This last moving average can be seen as

$$(I * k)(i) = \sum_{j=-n}^n I(i-j) \times K(j) = \frac{1}{N} \sum_{j=m}^{-m} I(i-j) \quad (3)$$

With $I(j)$ representing the value of the pixel at position j ,

$$K(j) = \begin{cases} \frac{1}{N} & \text{if } j \in \{-m, -m+1, \dots, 1, 0, 1, \dots, m-1, m\} \\ 0 & \text{else} \end{cases}$$

with $0 < m < n$.

This can be generalized into the 2D images

Left I and Right $I * K$

0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	90	0	0	0
0	0	0	90	90	90	90	90	90	0	0	0
0	0	0	90	90	90	90	90	90	0	0	0
0	0	0	90	0	90	90	90	90	0	0	0
0	0	0	90	90	90	90	90	90	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0

			0								

This can be generalized into the 2D images

Left I and Right $I * K$

0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	90	0	0
0	0	0	90	0	90	90	90	90	0	0
0	0	0	90	90	90	90	90	90	0	0
0	0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0

0	10									

This can be generalized into the 2D images

Left I and Right $I * K$

0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	90	90	0	0
0	0	0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0

0	10	20									

This can be generalized into the 2D images

Left I and Right $I * K$

0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0	0	0
0	0	0	90	90	90	90	90	0	0	0	0
0	0	0	90	90	90	90	90	0	0	0	0
0	0	0	90	0	90	90	90	0	0	0	0
0	0	0	90	90	90	90	90	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0

	0	10	20	30	30	30	20	10			
	0	20	40	60	60	60	40	20			
	0	30	60	90	90	90	60	30			
	0	30	50	80	80	90	60	30			
	0	30	50	80	80	90	60	30			
	0	20	30	50	50	60	40	20			
	10	20	30	30	30	30	20	10			
	10	10	10	0	0	0	0	0			

Moving average in 2D

Basically in 2D

- We can define different types of filter using the idea of weighted average

$$(I * K)(i, j) = \sum_{s=-m}^{-m} \sum_{l=-m}^m I(i - s, j - l) \times K(s, l) \quad (4)$$

Moving average in 2D

Basically in 2D

- We can define different types of filter using the idea of weighted average

$$(I * K)(i, j) = \sum_{s=-m}^{-m} \sum_{l=-m}^m I(i-s, j-l) \times K(s, l) \quad (4)$$

For example, the Box Filter

$$K = \frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} \text{"The Box Filter"} \quad (5)$$

Another Example

The Gaussian Filter

$$K = \begin{bmatrix} 0 & 1 & 2 & 1 & 0 \\ 1 & 3 & 4 & 3 & 1 \\ 2 & 5 & 9 & 5 & 2 \\ 1 & 3 & 5 & 3 & 1 \\ 0 & 1 & 2 & 1 & 0 \end{bmatrix}$$

Another Example

The Gaussian Filter

$$K = \begin{bmatrix} 0 & 1 & 2 & 1 & 0 \\ 1 & 3 & 4 & 3 & 1 \\ 2 & 5 & 9 & 5 & 2 \\ 1 & 3 & 5 & 3 & 1 \\ 0 & 1 & 2 & 1 & 0 \end{bmatrix}$$

Thus, we can define the concept of convolution

- Yes, using the previous ideas

Convolution

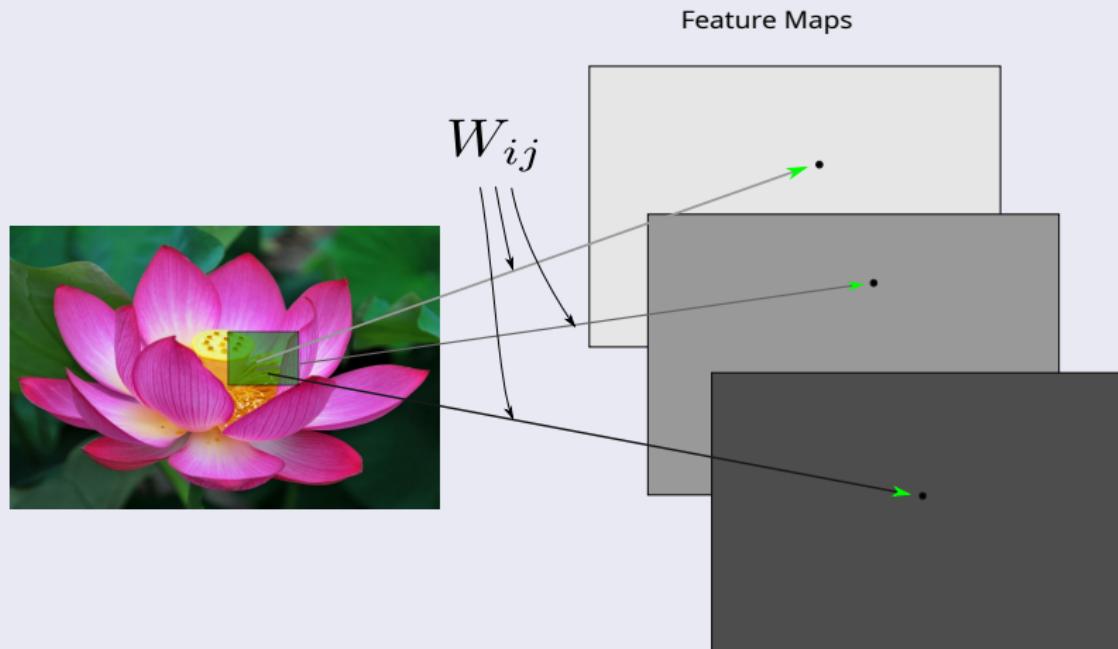
Definition

- Let $I : [a, b] \times [c, d] \rightarrow [0..255]$ be the image and $K : [e, f] \times [h, i] \rightarrow \mathbb{R}$ be the kernel. The output of Convolving I with K , denoted $I * K$ is

$$(I * K) [x, y] = \sum_{s=-n}^n \sum_{l=-n}^n I(x - s, y - l) \times K(s, l)$$

Now, why not to expand this idea

Imagine that a three channel image is splitted into a three feature map



Mathematically, we have the following

Map i

$$(I * k) [x, y, o] = \sum_{c=1}^3 \sum_{l=-n}^n \sum_{s=-n}^n I(x - l, y - s, c) \times k(l, s, c, o)$$

Mathematically, we have the following

Map i

$$(I * k) [x, y, o] = \sum_{c=1}^3 \sum_{l=-n}^n \sum_{s=-n}^n I(x - l, y - s, c) \times k(l, s, c, o)$$

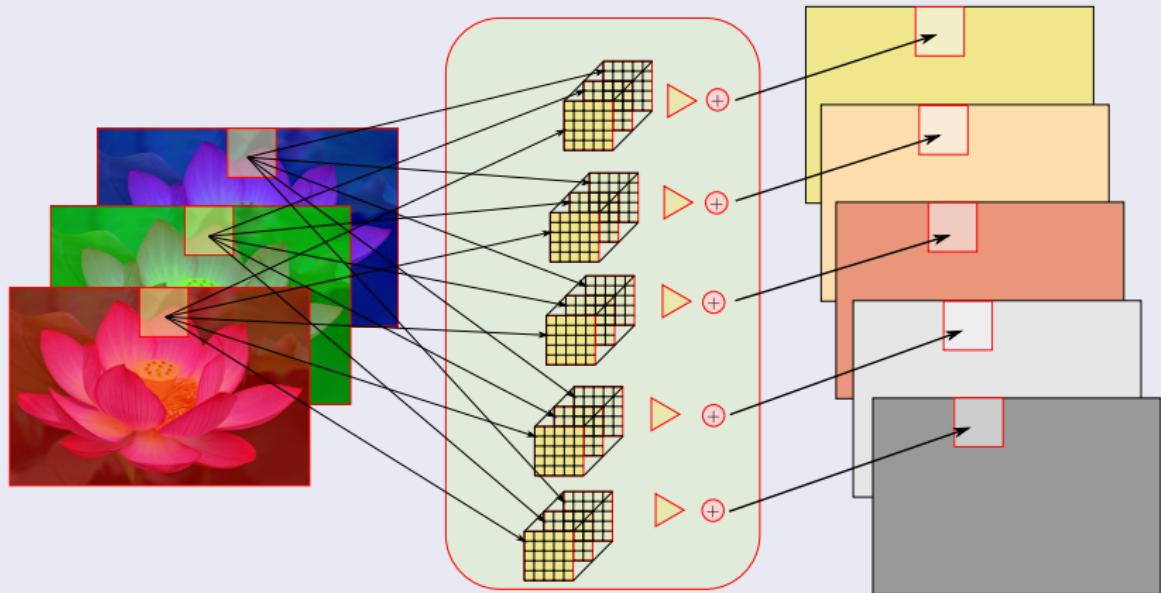
The convolution can be used as part of the follow architectures

- Filter
- Encoder
- Decoder
- etc

For Example, Encoder

We have the following situation

$$\sum_{c=1}^3 \sum_{l=-n}^n \sum_{s=-n}^n I(x - l, y - s, c) \times k(l, s, c, o)$$



Notation

We have the following

- $Y_j^{(l)}$ is a matrix representing the l layer and j^{th} feature map.
- $K_{ij}^{(l)}$ is the kernel filter with i^{th} kernel for layer j^{th} .

Notation

We have the following

- $Y_j^{(l)}$ is a matrix representing the l layer and j^{th} feature map.
- $K_{ij}^{(l)}$ is the kernel filter with i^{th} kernel for layer j^{th} .

Therefore

- We can see the Convolutional as a fusion of information from different feature maps.

$$\sum_{j=1}^{m_1^{(l-1)}} Y_j^{(l-1)} * K_{ij}^{(l)}$$

Thus, we have

Given a specific layer l , we have that i^{th} feature map in such layer equal to

$$Y_i^{(l)}(x, y) = B_i^{(l)}(x, y) + \sum_{j=1}^{m_1^{(l-1)}} \sum_{s=-ks}^{ks} \sum_{l=-ks}^{ks} Y_j^{(l-1)}(x-s, y-l) K_{ij}^{(l)}(x, y)$$

Thus, we have

Given a specific layer l , we have that i^{th} feature map in such layer equal to

$$Y_i^{(l)}(x, y) = B_i^{(l)}(x, y) + \sum_{j=1}^{m_1^{(l-1)}} \sum_{s=-ks}^{ks} \sum_{l=-ks}^{ks} Y_j^{(l-1)}(x-s, y-l) K_{ij}^{(l)}(x, y)$$

Where

- $Y_i^{(l)}$ is the i^{th} feature map in layer l .

Thus, we have

Given a specific layer l , we have that i^{th} feature map in such layer equal to

$$Y_i^{(l)}(x, y) = B_i^{(l)}(x, y) + \sum_{j=1}^{m_1^{(l-1)}} \sum_{s=-ks}^{ks} \sum_{l=-ks}^{ks} Y_j^{(l-1)}(x-s, y-l) K_{ij}^{(l)}(x, y)$$

Where

- $Y_i^{(l)}$ is the i^{th} feature map in layer l .
- $B_i^{(l)}$ is the bias matrix for output j .

Thus, we have

Given a specific layer l , we have that i^{th} feature map in such layer equal to

$$Y_i^{(l)}(x, y) = B_i^{(l)}(x, y) + \sum_{j=1}^{m_1^{(l-1)}} \sum_{s=-ks}^{ks} \sum_{l=-ks}^{ks} Y_j^{(l-1)}(x-s, y-l) K_{ij}^{(l)}(x, y)$$

Where

- $Y_i^{(l)}$ is the i^{th} feature map in layer l .
- $B_i^{(l)}$ is the bias matrix for output j .
- $K_{ij}^{(l)}$ is the filter of size $\left[2h_1^{(l)} + 1\right] \times \left[2h_2^{(l)} + 1\right]$.

Thus, we have

Given a specific layer l , we have that i^{th} feature map in such layer equal to

$$Y_i^{(l)}(x, y) = B_i^{(l)}(x, y) + \sum_{j=1}^{m_1^{(l-1)}} \sum_{s=-ks}^{ks} \sum_{l=-ks}^{ks} Y_j^{(l-1)}(x-s, y-l) K_{ij}^{(l)}(x, y)$$

Where

- $Y_i^{(l)}$ is the i^{th} feature map in layer l .
- $B_i^{(l)}$ is the bias matrix for output j .
- $K_{ij}^{(l)}$ is the filter of size $[2h_1^{(l)} + 1] \times [2h_2^{(l)} + 1]$.

Thus

- The input of layer l comprises $m_1^{(l-1)}$ feature maps from the previous layer, each of size $m_2^{(l-1)} \times m_3^{(l-1)}$

Therefore

The output of layer l

- It consists $m_1^{(l)}$ feature maps of size $m_2^{(l)} \times m_3^{(l)}$

Therefore

The output of layer l

- It consists $m_1^{(l)}$ feature maps of size $m_2^{(l)} \times m_3^{(l)}$

Something Notable

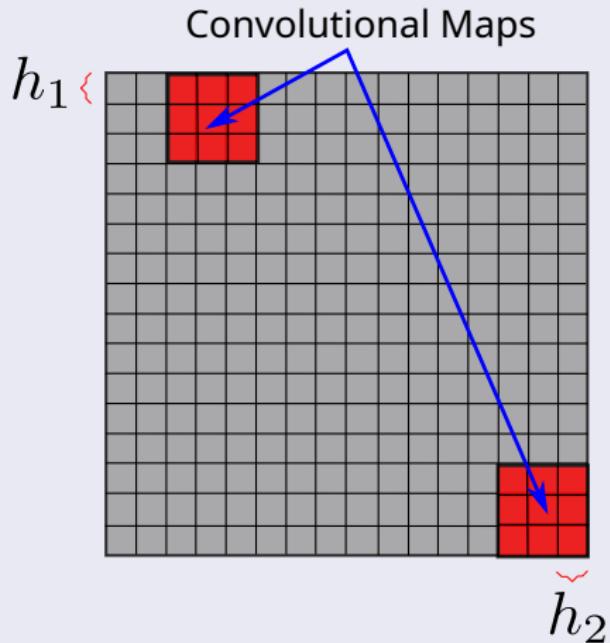
- $m_2^{(l)}$ and $m_3^{(l)}$ are influenced by border effects.
- Therefore, the output feature maps when the Convolutional sum is defined properly have size

$$m_2^{(l)} = m_2^{(l-1)} - 2h_1^{(l)}$$

$$m_3^{(l)} = m_3^{(l-1)} - 2h_2^{(l)}$$

Why? The Border

Example



Special Case

When $l = 1$

The input is a single image I consisting of one or more channels.

Thus

We have

Each feature map $Y_i^{(l)}$ in layer l consists of $m_1^{(l)} \cdot m_2^{(l)}$ units arranged in a two dimensional array.

Thus

We have

Each feature map $Y_i^{(l)}$ in layer l consists of $m_1^{(l)} \cdot m_2^{(l)}$ units arranged in a two dimensional array.

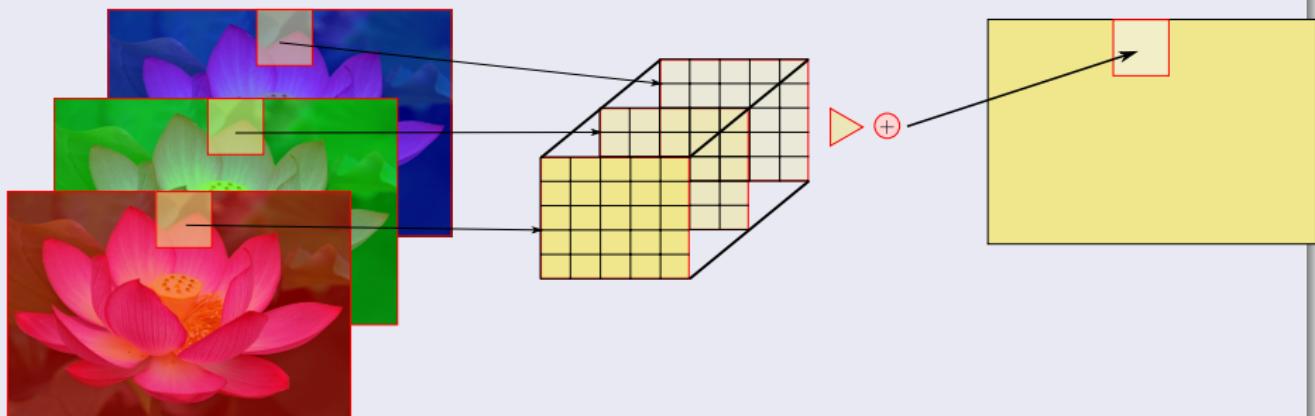
Thus, the unit at position (x, y) computes

$$\begin{aligned} (Y_i^{(l)})_{x,y} &= (B_i^{(l)})_{x,y} + \sum_{j=1}^{m_1^{(l-1)}} (K_{ij}^{(l)} * Y_j^{(l-1)})_{x,y} \\ &= (B_i^{(l)})_{x,y} + \sum_{j=1}^{m_1^{(l-1)}} \sum_{k=-h_1^{(l)}}^{h_1^{(l)}} \sum_{t=-h_2^{(l)}}^{h_2^{(l)}} (K_{ij}^{(l)})_{k,t} (Y_j^{(l-1)})_{x-k, x-t} \end{aligned}$$

This equation is basically

A use of the channels to generate a new channel

$$\left(B_i^{(l)} \right)_{x,y} + \sum_{j=1}^{m_1^{(l-1)}} \sum_{k=-h_1^{(l)}}^{h_1^{(l)}} \sum_{t=-h_2^{(l)}}^{h_2^{(l)}} \left(K_{ij}^{(l)} \right)_{k,t} \left(Y_j^{(l-1)} \right)_{x-k, x-t}$$



Given the fact that flipping the matrix is costly

We actually use Cross Correlation

$$(I \star K)(i, j) = \sum_{s=m}^{-m} \sum_{l=-m}^m I(i + s, j + l) \times K(s, l)$$

Given the fact that flipping the matrix is costly

We actually use Cross Correlation

$$(I \star K)(i, j) = \sum_{s=m}^{-m} \sum_{l=-m}^m I(i + s, j + l) \times K(s, l)$$

Actually given the adjustment during backpropagation

- We can have cross-correlation, convolution or whatever...

Now, we need to have an inverse operation to Convolution

For this

- We need to talk about the transpose of the Convolution.

Now, we need to have an inverse operation to Convolution

For this

- We need to talk about the transpose of the Convolution.

This is known as Deconvolution

- Which has a large history...

Here, an interesting case

Only a Historical Note

- The foundations for deconvolution came from Norbert Wiener of the Massachusetts Institute of Technology in his book “Extrapolation, Interpolation, and Smoothing of Stationary Time Series” (1949)

Here, an interesting case

Only a Historical Note

- The foundations for deconvolution came from Norbert Wiener of the Massachusetts Institute of Technology in his book “Extrapolation, Interpolation, and Smoothing of Stationary Time Series” (1949)

Basically, it tries to solve the following equation with $Y^{(l)}$ unknown layer that we want to recover

$$Y_i^{(l)} * K_{ij}^{(l)} = Y_j^{(l-1)}$$

Transposed Convolutions

This is the solution

- To go from a lower dimensional space to a higher dimensional space
 - ▶ While maintaining the connectivity pattern

Transposed Convolutions

This is the solution

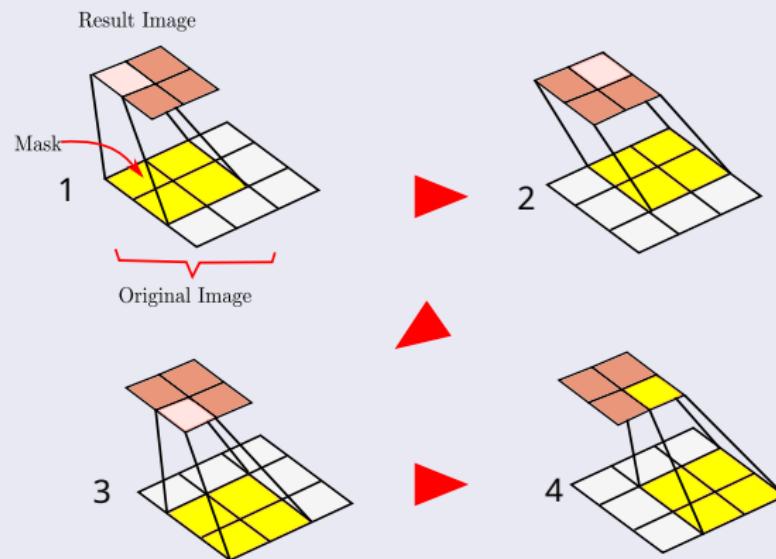
- To go from a lower dimensional space to a higher dimensional space
 - ▶ While maintaining the connectivity pattern

The solution

- Transposed convolutions – also called fractionally strided convolutions

Convolution as Matrix Computations

Imagine the convolution with no padding



We have the following idea

We unroll the image

$$\begin{pmatrix} i_{0,0} & i_{0,1} & i_{0,2} \\ i_{1,0} & i_{1,1} & i_{1,2} \\ i_{2,0} & i_{2,1} & i_{2,2} \end{pmatrix}$$

We have the following idea

We unroll the image

$$\begin{pmatrix} i_{0,0} & i_{0,1} & i_{0,2} \\ i_{1,0} & i_{1,1} & i_{1,2} \\ i_{2,0} & i_{2,1} & i_{2,2} \end{pmatrix}$$

We have the image as

$$\begin{pmatrix} i_{0,0} & i_{0,1} & i_{0,2} & i_{1,0} & i_{1,1} & i_{1,2} & i_{2,0} & i_{2,1} & i_{2,2} \end{pmatrix}^T$$

Now, the mask

The Mask can be arranged as

$$\begin{pmatrix} i_{0,0} & i_{0,1} & i_{0,2} & i_{1,0} & i_{1,1} & i_{1,2} & i_{2,0} & i_{2,1} & i_{2,2} \end{pmatrix}^T$$

$$\begin{pmatrix} w_{0,0} & w_{0,1} \\ w_{1,0} & w_{1,1} \end{pmatrix} \Rightarrow \begin{pmatrix} w_{0,0} & w_{0,1} & 0 & w_{1,0} & w_{1,1} & 0 & 0 & 0 & 0 \\ 0 & w_{0,0} & w_{0,1} & 0 & w_{1,0} & w_{1,1} & 0 & 0 & 0 \\ 0 & 0 & 0 & w_{0,0} & w_{0,1} & 0 & w_{1,0} & w_{1,1} & 0 \\ 0 & 0 & 0 & 0 & w_{0,0} & w_{0,1} & 0 & w_{1,0} & w_{1,1} \end{pmatrix}$$

Here

We have

- Images $m \times n$ get converted to flat vectors of $m * n$

Here

We have

- Images $m \times n$ get converted to flat vectors of $m * n$

Then using the Toplitz Matrix

$$W_{t \times p} I_f^{p=m*n} = I_f^t$$

Then the Transpose Matrix is simply that

THE TRANSPOSE!!! To reverse the original convolution

$$\begin{pmatrix} w'_{0,0} & 0 & 0 & 0 \\ w'_{0,1} & w'_{0,0} & 0 & 0 \\ 0 & w'_{0,1} & 0 & 0 \\ w'_{1,0} & 0 & w'_{0,0} & 0 \\ w'_{1,1} & w'_{1,0} & w'_{0,1} & w'_{0,0} \\ 0 & w'_{1,1} & 0 & w'_{0,1} \\ 0 & 0 & w'_{1,0} & 0 \\ 0 & 0 & w'_{1,1} & w'_{1,0} \\ 0 & 0 & 0 & w'_{1,1} \end{pmatrix}$$

Outline

1 Introduction

- The Long Path
- The Problem of Image Processing
- Multilayer Neural Network Classification
- Drawbacks
- Possible Solution

2 Convolutional Networks

- History
- Local Connectivity
- Sharing Parameters

3 Layers

- Convolutional Layer
- Convolutional Architectures
- A Little Bit of Notation

- The Reality on Models - The Use of Correlation
- Deconvolution Layer

• Non-Linearity Layer

- Fixing the Problem, ReLu function
- Back to the Non-Linearity Layer
- Rectification Layer
- Local Contrast Normalization Layer
- Sub-sampling and Pooling
- Strides
- Finally, The Fully Connected Layer

4 An Example of CNN

- The Proposed Architecture
- Backpropagation
- Deriving $w_{r,s,k}$
- Deriving the Kernel Filters

As in a Multilayer Perceptron

We use a non-linearity

- However, there is a drawback when using Back-Propagation under a sigmoid function

$$s(x) = \frac{1}{1 + e^{-x}}$$

As in a Multilayer Perceptron

We use a non-linearity

- However, there is a drawback when using Back-Propagation under a sigmoid function

$$s(x) = \frac{1}{1 + e^{-x}}$$

Because if we imagine a Convolutional Network as a series of layer functions f_i

$$y(A) = f_t \circ f_{t-1} \circ \cdots \circ f_2 \circ f_1(A)$$

With f_t is the last layer.

As in a Multilayer Perceptron

We use a non-linearity

- However, there is a drawback when using Back-Propagation under a sigmoid function

$$s(x) = \frac{1}{1 + e^{-x}}$$

Because if we imagine a Convolutional Network as a series of layer functions f_i

$$y(A) = f_t \circ f_{t-1} \circ \cdots \circ f_2 \circ f_1(A)$$

With f_t is the last layer.

Therefore, we finish with a sequence of derivatives

$$\frac{\partial y(A)}{\partial w_{1i}} = \frac{\partial f_t(f_{t-1})}{\partial f_{t-1}} \cdot \frac{\partial f_{t-1}(f_{t-2})}{\partial f_{t-2}} \cdot \dots \cdot \frac{\partial f_2(f_1)}{\partial f_2} \cdot \frac{\partial f_1(A)}{\partial w_{1i}}$$

Therefore

Given the commutativity of the product

- You could put together the derivative of the sigmoid's

$$f(x) = \frac{ds(x)}{dx} = \frac{e^{-x}}{(1 + e^{-x})^2}$$

Therefore

Given the commutativity of the product

- You could put together the derivative of the sigmoid's

$$f'(x) = \frac{ds(x)}{dx} = \frac{e^{-x}}{(1 + e^{-x})^2}$$

Therefore, deriving again

$$\frac{df(x)}{dx} = -\frac{e^{-x}}{(1 + e^{-x})^2} + \frac{2(e^{-x})^2}{(1 + e^{-x})^3}$$

Therefore

Given the commutativity of the product

- You could put together the derivative of the sigmoid's

$$f'(x) = \frac{ds(x)}{dx} = \frac{e^{-x}}{(1 + e^{-x})^2}$$

Therefore, deriving again

$$\frac{df(x)}{dx} = -\frac{e^{-x}}{(1 + e^{-x})^2} + \frac{2(e^{-x})^2}{(1 + e^{-x})^3}$$

After making $\frac{df(x)}{dx} = 0$

- We have the maximum is at $x = 0$

Therefore

The maximum for the derivative of the sigmoid

- $f'(0) = 0.25$

Therefore

The maximum for the derivative of the sigmoid

- $f'(0) = 0.25$

Therefore, Given a **Deep** Convolutional Network

- We could finish with

$$\lim_{k \rightarrow \infty} \left(\frac{ds(x)}{dx} \right)^k = \lim_{k \rightarrow \infty} (0.25)^k \rightarrow 0$$

Therefore

The maximum for the derivative of the sigmoid

- $f'(0) = 0.25$

Therefore, Given a **Deep** Convolutional Network

- We could finish with

$$\lim_{k \rightarrow \infty} \left(\frac{ds(x)}{dx} \right)^k = \lim_{k \rightarrow \infty} (0.25)^k \rightarrow 0$$

A vanishing derivative

- Making quite difficult to do train a deeper network using this activation function

Thus

The need to introduce a new function

$$f(x) = x^+ = \max(0, x)$$

Thus

The need to introduce a new function

$$f(x) = x^+ = \max(0, x)$$

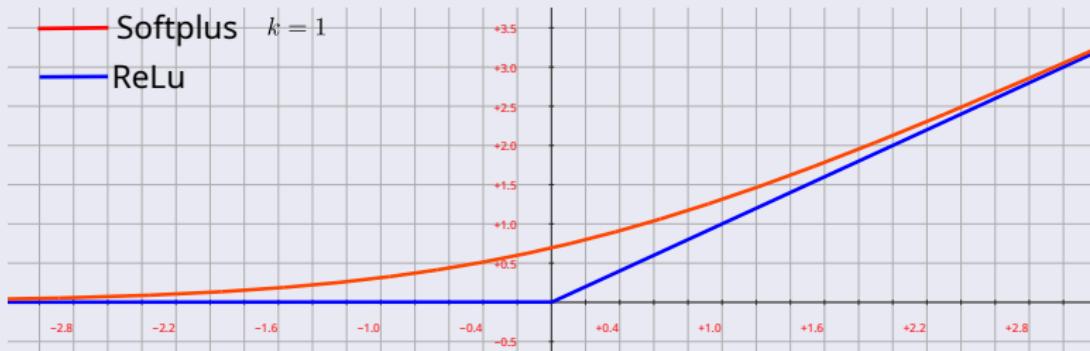
It is called ReLu or Rectifier

With a smooth approximation (Softplus function)

$$f(x) = \frac{\ln(1 + e^{kx})}{k}$$

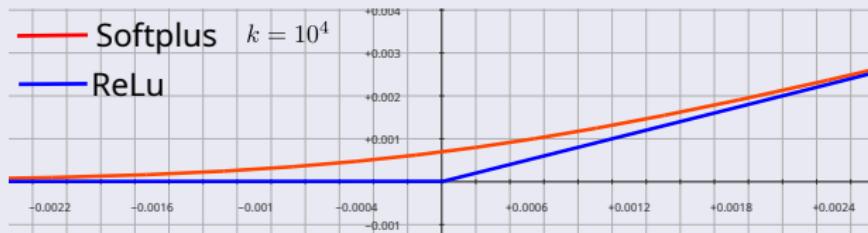
Therefore, we have

When $k = 1$



Increase k

When $k = 10^4$



Non-Linearity Layer

If layer l is a non-linearity layer

Its input is given by $m_1^{(l)}$ feature maps.

Non-Linearity Layer

If layer l is a non-linearity layer

Its input is given by $m_1^{(l)}$ feature maps.

What about the output

Its output comprises again $m_1^{(l)} = m_1^{(l-1)}$ feature maps

Non-Linearity Layer

If layer l is a non-linearity layer

Its input is given by $m_1^{(l)}$ feature maps.

What about the output

Its output comprises again $m_1^{(l)} = m_1^{(l-1)}$ feature maps

Each of them of size

$$m_2^{(l-1)} \times m_3^{(l-1)} \quad (6)$$

With $m_2^{(l)} = m_2^{(l-1)}$ and $m_3^{(l)} = m_3^{(l-1)}$.

Thus

With the final output

$$Y_i^{(l)} = f(Y_i^{(l-1)}) \quad (7)$$

Thus

With the final output

$$Y_i^{(l)} = f(Y_i^{(l-1)}) \quad (7)$$

Where

f is the activation function used in layer l and operates point wise.

Thus

With the final output

$$Y_i^{(l)} = f(Y_i^{(l-1)}) \quad (7)$$

Where

f is the activation function used in layer l and operates point wise.

You can also add a gain to compensate

$$Y_i^{(l)} = g_i f(Y_i^{(l-1)}) \quad (8)$$

Outline

1 Introduction

- The Long Path
- The Problem of Image Processing
- Multilayer Neural Network Classification
- Drawbacks
- Possible Solution

2 Convolutional Networks

- History
- Local Connectivity
- Sharing Parameters

3 Layers

- Convolutional Layer
- Convolutional Architectures
- A Little Bit of Notation

• The Reality on Models - The Use of Correlation

- Deconvolution Layer
- Non-Linearity Layer
- Fixing the Problem, ReLu function
- Back to the Non-Linearity Layer

• Rectification Layer

- Local Contrast Normalization Layer
- Sub-sampling and Pooling
- Strides
- Finally, The Fully Connected Layer

4 An Example of CNN

- The Proposed Architecture
- Backpropagation
- Deriving $w_{r,s,k}$
- Deriving the Kernel Filters

Rectification Layer, R_{abs}

Now a rectification layer

Then its input comprises $m_1^{(l)}$ feature maps of size $m_2^{(l-1)} \times m_3^{(l-1)}$.

Rectification Layer, R_{abs}

Now a rectification layer

Then its input comprises $m_1^{(l)}$ feature maps of size $m_2^{(l-1)} \times m_3^{(l-1)}$.

Then, the absolute value for each component of the feature maps is computed

$$Y_i^{(l)} = |Y_i^{(l)}| \quad (9)$$

Rectification Layer, R_{abs}

Now a rectification layer

Then its input comprises $m_1^{(l)}$ feature maps of size $m_2^{(l-1)} \times m_3^{(l-1)}$.

Then, the absolute value for each component of the feature maps is computed

$$Y_i^{(l)} = |Y_i^{(l)}| \quad (9)$$

Where the absolute value

It is computed point wise such that the output consists of $m_1^{(l)} = m_1^{(l-1)}$ feature maps unchanged in size.

Thus

We have that

Experiments show that rectification plays a central role in achieving good performance.

Thus

We have that

Experiments show that rectification plays a central role in achieving good performance.

You can find this in

K. Jarrett, K. Kavukcuogl, M. Ranzato, and Y. LeCun. What is the best multi-stage architecture for object recognition? In Computer Vision, International Conference on, pages 2146–2153, 2009.

Thus

We have that

Experiments show that rectification plays a central role in achieving good performance.

You can find this in

K. Jarrett, K. Kavukcuogl, M. Ranzato, and Y. LeCun. What is the best multi-stage architecture for object recognition? In Computer Vision, International Conference on, pages 2146–2153, 2009.

Remark

- Rectification could be included in the non-linearity layer.
- But also it can be seen as an independent layer.

Given that we are using Backpropagation

We need a soft approximation to $f(x) = |x|$

For this, we have

$$\frac{\partial f}{\partial x} = \text{sgn}(x)$$

- When $x \neq 0$. Why?

Given that we are using Backpropagation

We need a soft approximation to $f(x) = |x|$

For this, we have

$$\frac{\partial f}{\partial x} = \text{sgn}(x)$$

- When $x \neq 0$. Why?

We can use the following approximation

$$\text{sgn}(x) = 2 \left(\frac{\exp\{kx\}}{1 + \exp\{kx\}} \right) - 1$$

Given that we are using Backpropagation

We need a soft approximation to $f(x) = |x|$

For this, we have

$$\frac{\partial f}{\partial x} = \text{sgn}(x)$$

- When $x \neq 0$. Why?

We can use the following approximation

$$\text{sgn}(x) = 2 \left(\frac{\exp\{kx\}}{1 + \exp\{kx\}} \right) - 1$$

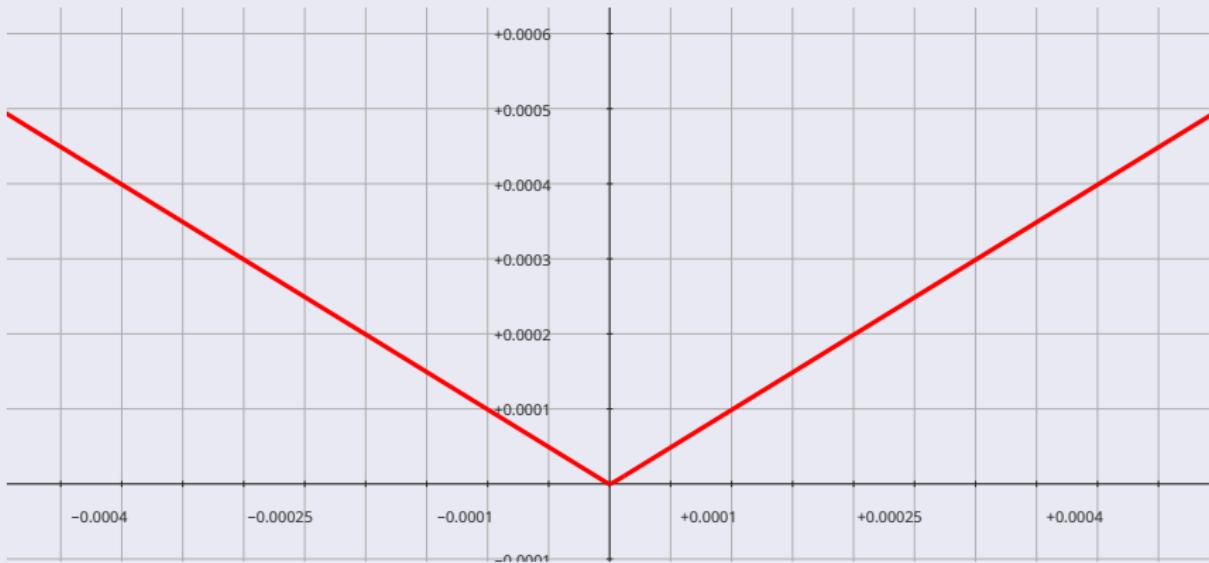
Therefore, we have by integration and working the C

$$f(x) = \frac{2}{k} \ln(1 + \exp\{kx\}) - x - \frac{2}{k} \ln(2)$$

We get the following situation

Something Notable

$$f(x) = \frac{2}{k} \ln(1 + \exp\{kx\}) - x - \frac{2}{k} \ln(2)$$



Outline

1 Introduction

- The Long Path
- The Problem of Image Processing
- Multilayer Neural Network Classification
- Drawbacks
- Possible Solution

2 Convolutional Networks

- History
- Local Connectivity
- Sharing Parameters

3 Layers

- Convolutional Layer
- Convolutional Architectures
- A Little Bit of Notation

• The Reality on Models - The Use of Correlation

- Deconvolution Layer
- Non-Linearity Layer
- Fixing the Problem, ReLu function
- Back to the Non-Linearity Layer
- Rectification Layer

• Local Contrast Normalization Layer

- Sub-sampling and Pooling
- Strides
- Finally, The Fully Connected Layer

4 An Example of CNN

- The Proposed Architecture
- Backpropagation
- Deriving $w_{r,s,k}$
- Deriving the Kernel Filters

Normalizing

Contrast normalization layer

The task of a local contrast normalization layer:

Normalizing

Contrast normalization layer

The task of a local contrast normalization layer:

- To enforce local competitiveness between adjacent units within a feature map.

Normalizing

Contrast normalization layer

The task of a local contrast normalization layer:

- To enforce local competitiveness between adjacent units within a feature map.
- To enforce competitiveness units at the same spatial location.

Normalizing

Contrast normalization layer

The task of a local contrast normalization layer:

- To enforce local competitiveness between adjacent units within a feature map.
- To enforce competitiveness units at the same spatial location.

We have two types of operations

- Subtractive Normalization.

Normalizing

Contrast normalization layer

The task of a local contrast normalization layer:

- To enforce local competitiveness between adjacent units within a feature map.
- To enforce competitiveness units at the same spatial location.

We have two types of operations

- Subtractive Normalization.
- Brightness Normalization.

Subtractive Normalization

Given $m_1^{(l-1)}$ feature maps of size $m_2^{(l-1)} \times m_3^{(l-1)}$

The output of layer l comprises $m_1^{(l)} = m_1^{(l-1)}$ feature maps unchanged in size.

Subtractive Normalization

Given $m_1^{(l-1)}$ feature maps of size $m_2^{(l-1)} \times m_3^{(l-1)}$

The output of layer l comprises $m_1^{(l)} = m_1^{(l-1)}$ feature maps unchanged in size.

With the operation

$$Y_i^{(l)} = Y_i^{(l-1)} - \sum_{j=1}^{m_1^{(l-1)}} K_{G(\sigma)} * Y_j^{(l-1)} \quad (10)$$

Subtractive Normalization

Given $m_1^{(l-1)}$ feature maps of size $m_2^{(l-1)} \times m_3^{(l-1)}$

The output of layer l comprises $m_1^{(l)} = m_1^{(l-1)}$ feature maps unchanged in size.

With the operation

$$Y_i^{(l)} = Y_i^{(l-1)} - \sum_{j=1}^{m_1^{(l-1)}} K_{G(\sigma)} * Y_j^{(l-1)} \quad (10)$$

With

$$\left(K_{G(\sigma)} \right)_{x,y} = \frac{1}{\sqrt{2\pi}\sigma^2} \exp \left\{ \frac{x^2 + y^2}{2\sigma^2} \right\} \quad (11)$$

Brightness Normalization

An alternative is to normalize the brightness in combination with the **rectified linear units**

$$\left(Y_i^{(l)} \right)_{x,y} = \frac{\left(Y_i^{(l-1)} \right)_{x,y}}{\left(\kappa + \lambda \sum_{j=1}^{m_1^{(l-1)}} \left(Y_j^{(l-1)} \right)_{x,y}^2 \right)^{\mu}} \quad (12)$$

Brightness Normalization

An alternative is to normalize the brightness in combination with the **rectified linear units**

$$\left(Y_i^{(l)} \right)_{x,y} = \frac{\left(Y_i^{(l-1)} \right)_{x,y}}{\left(\kappa + \lambda \sum_{j=1}^{m_1^{(l-1)}} \left(Y_j^{(l-1)} \right)_{x,y}^2 \right)^{\mu}} \quad (12)$$

Where

- κ, μ and λ are hyperparameters which can be set using a

$$f(x) = \frac{\ln(1 + e^{kx})}{k}$$

validation set.

Outline

1 Introduction

- The Long Path
- The Problem of Image Processing
- Multilayer Neural Network Classification
- Drawbacks
- Possible Solution

2 Convolutional Networks

- History
- Local Connectivity
- Sharing Parameters

3 Layers

- Convolutional Layer
- Convolutional Architectures
- A Little Bit of Notation

• The Reality on Models - The Use of Correlation

- Deconvolution Layer
- Non-Linearity Layer
- Fixing the Problem, ReLu function
- Back to the Non-Linearity Layer
- Rectification Layer
- Local Contrast Normalization Layer

• Sub-sampling and Pooling

- Strides
- Finally, The Fully Connected Layer

4 An Example of CNN

- The Proposed Architecture
- Backpropagation
- Deriving $w_{r,s,k}$
- Deriving the Kernel Filters

Sub-sampling Layer

Motivation

The motivation of subsampling the feature maps obtained by previous layers is robustness to noise and distortions.

Sub-sampling Layer

Motivation

The motivation of subsampling the feature maps obtained by previous layers is robustness to noise and distortions.

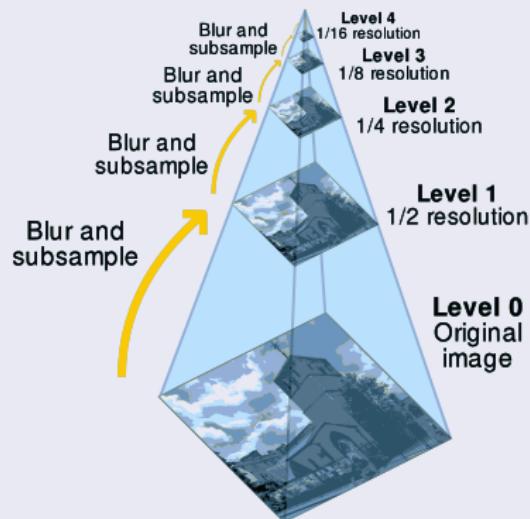
How?

- Normally, in traditional Convolutional Networks subsampling this is done by applying skipping factors!!!
- However, it is possible to combine subsampling with pooling and do it in a separate layer

Sub-sampling

The subsampling layer

- It seems to be acting as the well known sub-sampling pyramid



How is sub-sampling implemented?

We know that Image Pyramids

They were designed to find:

How is sub-sampling implemented?

We know that Image Pyramids

They were designed to find:

- ① filter-based representations to decompose images into information at multiple scales,

How is sub-sampling implemented?

We know that Image Pyramids

They were designed to find:

- ① filter-based representations to decompose images into information at multiple scales,
- ② To extract features/structures of interest,

How is sub-sampling implemented?

We know that Image Pyramids

They were designed to find:

- ① filter-based representations to decompose images into information at multiple scales,
- ② To extract features/structures of interest,
- ③ To attenuate noise.

How is sub-sampling implemented?

We know that Image Pyramids

They were designed to find:

- ① filter-based representations to decompose images into information at multiple scales,
- ② To extract features/structures of interest,
- ③ To attenuate noise.

Example of usage of this filters

- The SURF and SIFT filters

There are also other ways of doing this

subsampling can be done using so called skipping factors

$$s_1^{(l)} \text{ and } s_2^{(l)}$$

There are also other ways of doing this

subsampling can be done using so called skipping factors

$$s_1^{(l)} \text{ and } s_2^{(l)}$$

The basic idea is to skip a fixed number of pixels

Therefore the size of the output feature map is given by

$$m_2^{(l)} = \frac{m_2^{(l-1)} - 2h_1^{(l)}}{s_1^{(l)} + 1} \text{ and } m_3^{(l)} = \frac{m_3^{(l-1)} - 2h_2^{(l)}}{s_2^{(l)} + 1}$$

What is Pooling?

Pooling

- **Spatial pooling is way to compute image representation based on encoded local features.**

Pooling

Let l be a pooling layer

- It outputs from $m_i^{(l)} > m_i^{(l-1)}$ feature maps of reduced size.

Pooling

Let l be a pooling layer

- It outputs from $m_i^{(l)} > m_i^{(l-1)}$ feature maps of reduced size.

Pooling Operation

It operates by placing windows at non-overlapping positions in each feature map and keeping one value per window such that the feature maps are sub-sampled.

Thus

In the previous example

- All feature maps are pooled and sub-sampled individually.

Thus

In the previous example

- All feature maps are pooled and sub-sampled individually.

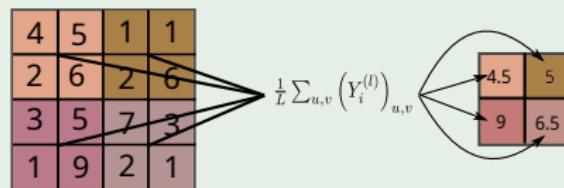
Each unit

- In one of the $m_1^{(l)} = 4$ output feature maps represents the average or the maximum within a fixed window of the corresponding feature map in layer $(l - 1)$.

Examples of pooling

Average pooling

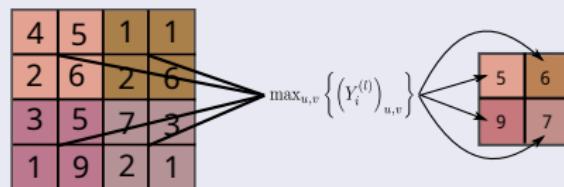
When using a boxcar filter, the operation is called average pooling and the layer denoted by P_A .



Examples of pooling

Max pooling

For max pooling, the maximum value of each window is taken. The layer is denoted by P_M .



An interesting property

Something notable depending in the pooling area

- “In all cases, pooling helps to make the representation become approximately invariant to small translations of the input. Invariance to translation means that if we translate the input by a small amount, the values of most of the pooled outputs do not change.”
 - ▶ Page 342, Ian Goodfellow, Introduction to Deep Learning, 2016 [8].

An interesting property

Something notable depending in the pooling area

- “In all cases, pooling helps to make the representation become approximately invariant to small translations of the input. Invariance to translation means that if we translate the input by a small amount, the values of most of the pooled outputs do not change.”
 - ▶ Page 342, Ian Goodfellow, Introduction to Deep Learning, 2016 [8].

The small amount

- In the case of the previous examples, **1 pixel**

Other Poolings

There are other types of pooling

- L_2 norm of a rectangular neighborhood
- Weighted average based on the distance from the central pixel

Other Poolings

There are other types of pooling

- L_2 norm of a rectangular neighborhood
- Weighted average based on the distance from the central pixel

However, we have another way of doing pooling

- Striding!!!

They started talking about substituting maxpooling for something called a Stride on the Convolution

$$(Y_i^{(l)})_{x,y} = (B_i^{(l)})_{x,y} + \sum_{j=1}^{m_1^{(l-1)}} \sum_{k=-h_1^{(l)}}^{h_1^{(l)}} \sum_{t=-h_2^{(l)}}^{h_2^{(l)}} (K_{ij}^{(l)})_{k,t} (Y_j^{(l-1)})_{x-k,x-t}$$

They started talking about substituting maxpooling for something called a Stride on the Convolution

$$(Y_i^{(l)})_{x,y} = (B_i^{(l)})_{x,y} + \sum_{j=1}^{m_1^{(l-1)}} \sum_{k=-h_1^{(l)}}^{h_1^{(l)}} \sum_{t=-h_2^{(l)}}^{h_2^{(l)}} (K_{ij}^{(l)})_{k,t} (Y_j^{(l-1)})_{x-k,x-t}$$

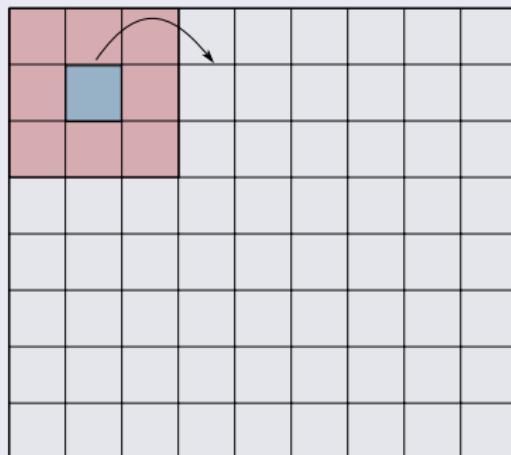
This is a Heuristic ...

- Basically you jump around by a factor r and t for the width and height of the layer
 - ▶ It was proposed to decrease memory usage...

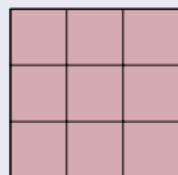
Example

Horizontal Stride

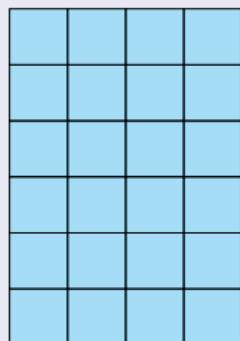
Horizontal Stride $r = 2$



*



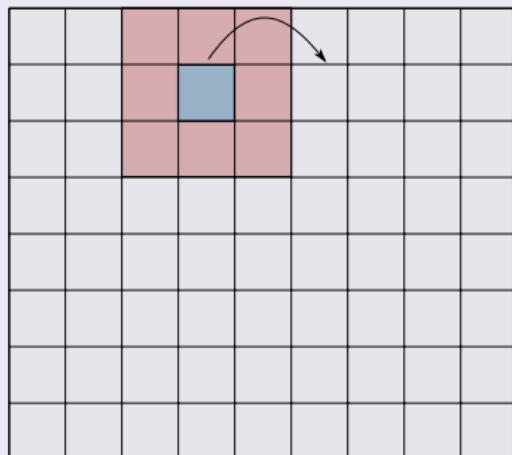
=



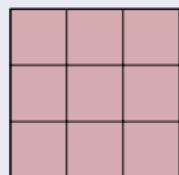
Example

Horizontal Stride

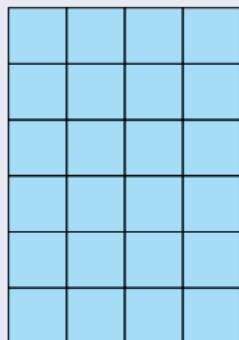
Horizontal Stride $r = 2$



*



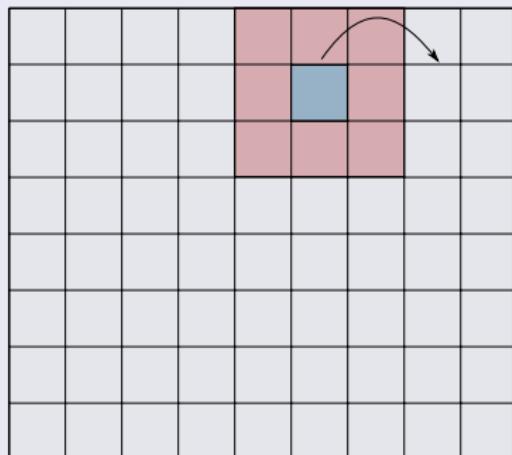
=



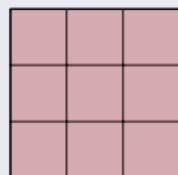
Example

Horizontal Stride

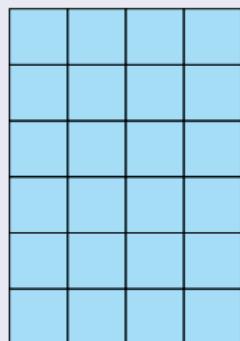
Horizontal Stride $r = 2$



*



=



There are attempts to understand its effects

At Convolution Level and using Tensors [10]

- “Take it in your stride: Do we need striding in CNNs?” by Chen Kong, Simon Lucey [11]

There are attempts to understand its effects

At Convolution Level and using Tensors [10]

- “Take it in your stride: Do we need striding in CNNs?” by Chen Kong, Simon Lucey [11]

Please read Kolda's Paper before you get into the other

- You need a little bit of notation...

Outline

1 Introduction

- The Long Path
- The Problem of Image Processing
- Multilayer Neural Network Classification
- Drawbacks
- Possible Solution

2 Convolutional Networks

- History
- Local Connectivity
- Sharing Parameters

3 Layers

- Convolutional Layer
- Convolutional Architectures
- A Little Bit of Notation

● The Reality on Models - The Use of Correlation

- Deconvolution Layer
- Non-Linearity Layer
- Fixing the Problem, ReLu function
- Back to the Non-Linearity Layer
- Rectification Layer
- Local Contrast Normalization Layer
- Sub-sampling and Pooling
- Strides

● Finally, The Fully Connected Layer

4 An Example of CNN

- The Proposed Architecture
- Backpropagation
- Deriving $w_{r,s,k}$
- Deriving the Kernel Filters

Actually

All the previous layers can be seen

- As the generation of an embedding
 - ▶ With Rotation and Centering done

Actually

All the previous layers can be seen

- As the generation of an embedding
 - ▶ With Rotation and Centering done

We still need a classifier

- A Classic selection, a MLP!!!

Fully Connected Layer

If a layer l is a fully connected layer

- If layer $(l - 1)$ is a fully connected layer, use the equation to compute the output of i^{th} unit at layer l :

$$z_i^{(l)} = \sum_{k=0}^{m^{(l)}} w_{i,k}^{(l)} y_k^{(l)} \text{ thus } y_i^{(l)} = f(z_i^{(l)})$$

Fully Connected Layer

If a layer l is a fully connected layer

- If layer $(l - 1)$ is a fully connected layer, use the equation to compute the output of i^{th} unit at layer l :

$$z_i^{(l)} = \sum_{k=0}^{m^{(l)}} w_{i,k}^{(l)} y_k^{(l)} \text{ thus } y_i^{(l)} = f(z_i^{(l)})$$

Otherwise

- Layer l expects $m_1^{(l-1)}$ feature maps of size $m_2^{(l-1)} \times m_3^{(l-1)}$ as input.

Then

Thus, the i^{th} unit in layer l computes

$$y_i^{(l)} = f(z_i^{(l)})$$

$$z_i^{(l)} = \sum_{j=1}^{m_1^{(l-1)}} \sum_{r=1}^{m_2^{(l-1)}} \sum_{s=1}^{m_3^{(l-1)}} w_{i,j,r,s}^{(l)} (Y_j^{(l-1)})_{r,s}$$

Here

Where $w_{i,j,r,s}^{(l)}$

- It denotes the weight connecting the unit at position (r, s) in the j^{th} feature map of layer $(l - 1)$ and the i^{th} unit in layer l .

Here

Where $w_{i,j,r,s}^{(l)}$

- It denotes the weight connecting the unit at position (r, s) in the j^{th} feature map of layer $(l - 1)$ and the i^{th} unit in layer l .

Something Notable

- In practice, Convolutional Layers are used to learn a feature hierarchy and one or more fully connected layers are used for classification purposes based on the computed features.

Basically

We can use a loss function at the output of such layer

$$L(\mathbf{W}) = \sum_{n=1}^N E_n(\mathbf{W}) = \sum_{n=1}^N \sum_{k=1}^K \left(y_{nk}^{(l)} - t_{nk} \right)^2 \text{ (Sum of Squared Error)}$$

$$L(\mathbf{W}) = \sum_{n=1}^N E_n(\mathbf{W}) = \sum_{n=1}^N \sum_{k=1}^K t_{nk} \log \left(y_{nk}^{(l)} \right) \text{ (Cross-Entropy Error)}$$

Basically

We can use a loss function at the output of such layer

$$L(\mathbf{W}) = \sum_{n=1}^N E_n(\mathbf{W}) = \sum_{n=1}^N \sum_{k=1}^K (y_{nk}^{(l)} - t_{nk})^2 \quad (\text{Sum of Squared Error})$$

$$L(\mathbf{W}) = \sum_{n=1}^N E_n(\mathbf{W}) = \sum_{n=1}^N \sum_{k=1}^K t_{nk} \log(y_{nk}^{(l)}) \quad (\text{Cross-Entropy Error})$$

Assuming \mathbf{W} the tensor used to represent all the possible weights

- We can use the Backpropagation idea as long we can apply the corresponding derivatives.

About this

As part of the seminar

- We are preparing a series of slides about Loss Functions...

Outline

1 Introduction

- The Long Path
- The Problem of Image Processing
- Multilayer Neural Network Classification
- Drawbacks
- Possible Solution

2 Convolutional Networks

- History
- Local Connectivity
- Sharing Parameters

3 Layers

- Convolutional Layer
- Convolutional Architectures
- A Little Bit of Notation

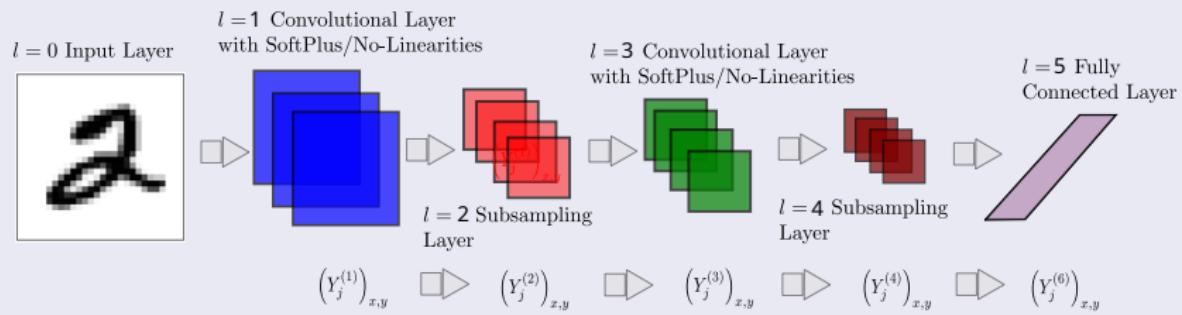
- The Reality on Models - The Use of Correlation
- Deconvolution Layer
- Non-Linearity Layer
- Fixing the Problem, ReLu function
- Back to the Non-Linearity Layer
- Rectification Layer
- Local Contrast Normalization Layer
- Sub-sampling and Pooling
- Strides
- Finally, The Fully Connected Layer

4 An Example of CNN

- The Proposed Architecture
- Backpropagation
- Deriving $w_{r,s,k}$
- Deriving the Kernel Filters

We have the following Architecture

Simplified Architecture by Jean LeCun “Backpropagation applied to handwritten zip code recognition”



Forward

We as always start at the begining

- The input is passed to the layer after a basic pre-processing as normalize $0, 1, 2, 3, 4, \dots 2^8 - 1$ values to values in the interval $[0, 1]$

Forward

We as always start at the begining

- The input is passed to the layer after a basic pre-processing as normalize $0, 1, 2, 3, 4, \dots 2^8 - 1$ values to values in the interval $[0, 1]$

Then we have

- The following forward

Therefore, we have

Layer $l = 1$

- This Layer is using a ReLu to output f with 3 channels

$$\left(Y_1^{(1)} \right)_{x,y} = \text{ReLU} \left[\left(B_1^{(l)} \right)_{x,y} + \sum_{k=-h_1^{(l)}}^{h_1^{(l)}} \sum_{t=-h_2^{(l)}}^{h_2^{(l)}} \left(K_{11}^{(l)} \right)_{k,t} \left(Y_1^{(l-1)} \right)_{x-k, x-t} \right]$$

$$\left(Y_2^{(1)} \right)_{x,y} = \text{ReLU} \left[\left(B_2^{(l)} \right)_{x,y} + \sum_{k=-h_1^{(l)}}^{h_1^{(l)}} \sum_{t=-h_2^{(l)}}^{h_2^{(l)}} \left(K_{21}^{(l)} \right)_{k,t} \left(Y_1^{(l-1)} \right)_{x-k, x-t} \right]$$

$$\left(Y_3^{(1)} \right)_{x,y} = \text{ReLU} \left[\left(B_3^{(l)} \right)_{x,y} + \sum_{k=-h_1^{(l)}}^{h_1^{(l)}} \sum_{t=-h_2^{(l)}}^{h_2^{(l)}} \left(K_{31}^{(l)} \right)_{k,t} \left(Y_1^{(l-1)} \right)_{x-k, x-t} \right]$$

Layer $l = 2$

We have a maxpooling of size 2×2

$$\left(Y_i^{(2)} \right)_{x',y'} = \max \left\{ \left(Y_i^{(l-1)} \right)_{x,y}, \left(Y_i^{(l-1)} \right)_{x+1,y}, \left(Y_i^{(l-1)} \right)_{x,y+1}, \left(Y_i^{(l-1)} \right)_{x+1,y+1} \right\}$$

Then, you repeat the previous process

Thus we obtain a reduced convoluted version $Y_m^{(3)}$ of the $Y_n^{(4)}$ convolution and maxpooling

- Thus, we use those as inputs for the fully connected layer of input.

The fully connected layer

Now assuming a single $k = 1$ neuron

$$z_1^{(5)} = \sum_{k=1}^9 \sum_{r=1}^{m_2^{(6)}} \sum_{s=1}^{m_3^{(6)}} w_{r,s,k}^{(5)} (Y_k^{(4)})_{r,s}$$
$$y_1^{(6)} = f(z_1^{(5)})$$

We have for simplicity sake

That our final cost function is equal to

$$L = \frac{1}{2} \left(y_1^{(6)} - t_1^{(6)} \right)^2$$

Outline

1 Introduction

- The Long Path
- The Problem of Image Processing
- Multilayer Neural Network Classification
- Drawbacks
- Possible Solution

2 Convolutional Networks

- History
- Local Connectivity
- Sharing Parameters

3 Layers

- Convolutional Layer
- Convolutional Architectures
- A Little Bit of Notation

- The Reality on Models - The Use of Correlation
- Deconvolution Layer
- Non-Linearity Layer
- Fixing the Problem, ReLu function
- Back to the Non-Linearity Layer
- Rectification Layer
- Local Contrast Normalization Layer
- Sub-sampling and Pooling
- Strides
- Finally, The Fully Connected Layer

4 An Example of CNN

- The Proposed Architecture
- **Backpropagation**
 - Deriving $w_{r,s,k}$
 - Deriving the Kernel Filters

After collecting all input/output

Therefore

- We have using sum of squared errors (loss function):

$$L = \frac{1}{2} \left(y_1^{(6)} - t_1^{(6)} \right)^2$$

After collecting all input/output

Therefore

- We have using sum of squared errors (loss function):

$$L = \frac{1}{2} \left(y_1^{(6)} - t_1^{(6)} \right)^2$$

Therefore, we can obtain

$$\frac{\partial L}{\partial w_{r,s,k}^{(5)}} = \frac{1}{2} \times \frac{\partial \left(y_1^{(6)} - t_1^{(6)} \right)^2}{\partial w_{r,s,k}^{(5)}}$$

Therefore

We get in the first part of the equation

$$\frac{\partial \left(y_1^{(6)} - t_1^{(6)} \right)^2}{\partial w_{r,s,k}^{(5)}} = \left(y_1^{(6)} - t_1^{(6)} \right) \frac{\partial y_1^{(6)}}{\partial w_{r,s,k}^{(5)}}$$

Therefore

We get in the first part of the equation

$$\frac{\partial \left(y_1^{(6)} - t_1^{(6)} \right)^2}{\partial w_{r,s,k}^{(5)}} = \left(y_1^{(6)} - t_1^{(6)} \right) \frac{\partial y_1^{(6)}}{\partial w_{r,s,k}^{(5)}}$$

With

$$y_1^{(6)} = ReLu \left(z_1^{(5)} \right)$$

Therefore

We have

$$\frac{\partial y_1^{(6)}}{\partial w_{r,s,k}^{(5)}} = \frac{\partial ReLu(z_1^{(5)})}{\partial z_1^{(5)}} \times \frac{\partial z_1^{(5)}}{\partial w_{r,s,k}^{(5)}}$$

Therefore

We have

$$\frac{\partial y_1^{(6)}}{\partial w_{r,s,k}^{(5)}} = \frac{\partial ReLu(z_1^{(5)})}{\partial z_1^{(5)}} \times \frac{\partial z_1^{(5)}}{\partial w_{r,s,k}^{(5)}}$$

Therefore if we use the approximation

$$\frac{\partial f(z_1^{(5)})}{\partial z_1^{(5)}} = \frac{e^{kz_1^{(5)}}}{(1 + e^{kz_1^{(5)}})}$$

Now, we need to derive $\frac{\partial z_1^{(5)}}{\partial w_{r,s,k}^{(5)}}$

We know that

$$z_1^{(5)} = \sum_{k=1}^9 \sum_{r=1}^{m_2^{(6)}} \sum_{s=1}^{m_3^{(6)}} w_{r,s,k}^{(5)} \left(Y_k^{(4)} \right)_{r,s}$$

Now, we need to derive $\frac{\partial z_1^{(5)}}{\partial w_{r,s,k}^{(5)}}$

We know that

$$z_1^{(5)} = \sum_{k=1}^9 \sum_{r=1}^{m_2^{(6)}} \sum_{s=1}^{m_3^{(6)}} w_{r,s,k}^{(5)} \left(Y_k^{(4)} \right)_{r,s}$$

Now, we need to derive $\frac{\partial z_1^{(5)}}{\partial w_{r,s,k}^{(5)}}$

We know that

$$z_1^{(5)} = \sum_{k=1}^9 \sum_{r=1}^{m_2^{(6)}} \sum_{s=1}^{m_3^{(6)}} w_{r,s,k}^{(5)} \left(Y_k^{(4)} \right)_{r,s}$$

Finally

$$\frac{\partial z_1^{(5)}}{\partial w_{r,s,k}^{(5)}} = \left(Y_k^{(4)} \right)_{r,s}$$

Maxpooling

This is not derived after all, but we go directly for the max term

- Assume you get the max element for $f = 1, 2, \dots, 9$ and $j = 1$

$$\left(Y_f^{(3)} \right)_{x,y} = \left(B_f^{(3)} \right)_{x,y} + \sum_{k=-h_1^{(l)}}^{h_1^{(l)}} \sum_{t=-h_2^{(l)}}^{h_2^{(l)}} \left(K_{f1}^{(3)} \right)_{k,t} \left(Y_1^{(2)} \right)_{x-k, y-t}$$

Therefore

We have then

$$\frac{\partial L}{\partial \left(K_{f1}^{(3)} \right)_{k,t}} = \frac{1}{2} \times \frac{\partial \left(y_1^{(6)} - t_1^{(6)} \right)^2}{\partial \left(K_{f1}^{(3)} \right)_{k,t}}$$

Therefore

We have then

$$\frac{\partial L}{\partial \left(K_{f1}^{(3)} \right)_{k,t}} = \frac{1}{2} \times \frac{\partial \left(y_1^{(6)} - t_1^{(6)} \right)^2}{\partial \left(K_{f1}^{(3)} \right)_{k,t}}$$

We have the following chain of derivations given

$$\left(Y_f^{(4)} \right)_{x,y} = f \left[\left(Y_f^{(3)} \right)_{x,y} \right]$$

$$\frac{\partial L}{\partial \left(K_{f1}^{(3)} \right)_{k,t}} = \left(y_1^{(6)} - t_1^{(6)} \right) \frac{\partial f \left(z_1^{(5)} \right)}{\partial z_1^{(5)}} \times \sum_{x,y \in Y_f^{(4)}} \left[\frac{\partial z_i^{(5)}}{\partial \left(Y_f^{(4)} \right)_{x,y}} \times \frac{\partial f \left[\left(Y_f^{(3)} \right)_{x,y} \right]}{\partial \left(K_{f1}^{(3)} \right)_{k,t}} \right]$$

Therefore

We have

$$\frac{\partial z_i^{(5)}}{\partial \left(Y_f^{(3)}\right)_{x,y}} = w_{x,y,f}^{(5)}$$

Therefore

We have

$$\frac{\partial z_i^{(5)}}{\partial \left(Y_f^{(3)}\right)_{x,y}} = w_{x,y,f}^{(5)}$$

Then assuming that

$$\left(Y_f^{(3)}\right)_{x,y} = \left(B_f^{(3)}\right)_{x,y} + \sum_{k=-h_1^{(l)}}^{h_1^{(l)}} \sum_{t=-h_2^{(l)}}^{h_2^{(l)}} \left(K_{f1}^{(3)}\right)_{k,t} \left(Y_1^{(2)}\right)_{x-k,x-t}$$

Therefore

We have

$$\frac{\partial f \left[\left(Y_f^{(3)} \right)_{x,y} \right]}{\partial \left(K_{f1}^{(3)} \right)_{k,t}} = \frac{\partial f \left[\left(Y_f^{(3)} \right)_{x,y} \right]}{\partial \left(Y_f^{(3)} \right)_{x,y}} \times \frac{\partial \left(Y_f^{(3)} \right)_{x,y}}{\partial \left(K_{f1}^{(3)} \right)_{k,t}}$$

Therefore

We have

$$\frac{\partial f \left[\left(Y_f^{(3)} \right)_{x,y} \right]}{\partial \left(K_{f1}^{(3)} \right)_{k,t}} = \frac{\partial f \left[\left(Y_f^{(3)} \right)_{x,y} \right]}{\partial \left(Y_f^{(3)} \right)_{x,y}} \times \frac{\partial \left(Y_f^{(3)} \right)_{x,y}}{\partial \left(K_{f1}^{(3)} \right)_{k,t}}$$

Then

$$\frac{\partial f \left[\left(Y_f^{(3)} \right)_{x,y} \right]}{\partial \left(Y_f^{(3)} \right)_{x,y}} = f' \left[\left(Y_f^{(3)} \right)_{x,y} \right]$$

Finally, we have

The equation

$$\frac{\partial \left(Y_f^{(3)} \right)_{x,y}}{\partial \left(K_{f1}^{(3)} \right)_{k,t}} = \left(Y_1^{(2)} \right)_{x-k,x-t}$$

The Total Equation is

The equation

$$\frac{\partial L}{\partial \left(K_{f1}^{(3)} \right)_{k,t}} = \left(y_1^{(6)} - t_1^{(6)} \right) \underbrace{\frac{\partial f \left(z_1^{(5)} \right)}{\partial z_1^{(5)}}}_{\text{Relu Derivative}} \times \sum_{x,y \in Y_f^{(4)}} \left[w_{x,y,f}^{(5)} \times f' \left[\left(Y_f^{(3)} \right)_{x,y} \right] \times \left(Y_1^{(2)} \right)_{x,y} \right]$$

We have the following graphically at Convolution Layer

At Forward

Forward



0	0	1	1	0
1	0	1	1	1
0	0	1	1	1
0	0	0	0	1
1	1	0	0	0

$$Y_j^{(l-1)}$$

1	2	1
1	1	1
1	2	1

$$K^{(l-1)}$$

$$+ f()$$



4	8	10
3	5	6
4	4	6

$$Y_j^{(l)}$$

We have the following graphically at the Convolution Layer

At Forward

Backward

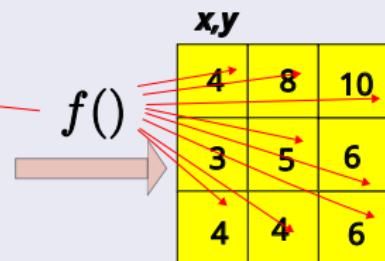


0	0	1	1	0
1	0	1	1	1
0	0	1	1	1
0	0	0	0	1
1	1	0	0	0

$Y_j^{(l-1)}$

h,t		
1	2	1
1	1	1
1	2	1

$K^{(l-1)}$



$Y_j^{(l)}$

The Other Equations

I will leave you to devise them

- They are a repetitive procedure.

The Other Equations

I will leave you to devise them

- They are a repetitive procedure.

The interesting case the average pooling

- The others are the stride and the deconvolution

- [1] A. Khan, A. Sohail, U. Zahoor, and A. S. Qureshi, "A survey of the recent architectures of deep convolutional neural networks," *Artificial Intelligence Review*, vol. 53, no. 8, pp. 5455–5516, 2020.
- [2] R. Szeliski, *Computer Vision: Algorithms and Applications*. Berlin, Heidelberg: Springer-Verlag, 1st ed., 2010.
- [3] S. Haykin, *Neural Networks and Learning Machines*. No. v. 10 in Neural networks and learning machines, Prentice Hall, 2009.
- [4] D. H. Hubel and T. N. Wiesel, "Receptive fields, binocular interaction and functional architecture in the cat's visual cortex," *The Journal of physiology*, vol. 160, no. 1, p. 106, 1962.
- [5] Y. LeCun, L. Bottou, Y. Bengio, P. Haffner, *et al.*, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.

- [6] W. Zhang, K. Itoh, J. Tanida, and Y. Ichioka, "Parallel distributed processing model with local space-invariant interconnections and its optical architecture," *Appl. Opt.*, vol. 29, pp. 4790–4797, Nov 1990.
- [7] J. J. Weng, N. Ahuja, and T. S. Huang, "Learning recognition and segmentation of 3-d objects from 2-d images," in *1993 (4th) International Conference on Computer Vision*, pp. 121–128, IEEE, 1993.
- [8] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. The MIT Press, 2016.
- [9] J. T. Springenberg, A. Dosovitskiy, T. Brox, and M. Riedmiller, "Striving for simplicity: The all convolutional net," 2015.
- [10] T. G. Kolda and B. W. Bader, "Tensor decompositions and applications," *SIAM review*, vol. 51, no. 3, pp. 455–500, 2009.
- [11] C. Kong and S. Lucey, "Take it in your stride: Do we need striding in cnns?," *arXiv preprint arXiv:1712.02502*, 2017.