

Introduction to Neural Networks and Deep Learning

Introduction to the Convolutional Network

Andres Mendez-Vazquez

March 28, 2021

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
 - Deconvolution Layer
 - Alternating Minimization
- 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
- Normalization Layer AKA Batch Normalization
- 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
 - Deconvolution Layer
 - Alternating Minimization
- 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
- Normalization Layer AKA Batch Normalization
- Finally, The Fully Connected Layer

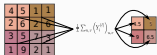
4 An Example of CNN

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

The Long Path [1]

A Small History of a Revolution

$$Y_i^{(l)}(x, y) = B_i^{(l)}(x, y) + \sum_{j=1}^{n_i-1} \sum_{k=-s_j}^{s_j} \sum_{l=-b_j}^{b_j} Y_j^{(l-1)}(x-s, y-l) K_{ij}^{(l)}(x, y)$$



Beyond

Complex Architectures and
The Attention Revolution

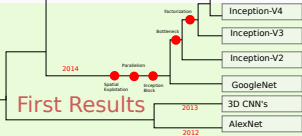
Residual and Multipath
Architectures

The Beginnig of Attention?

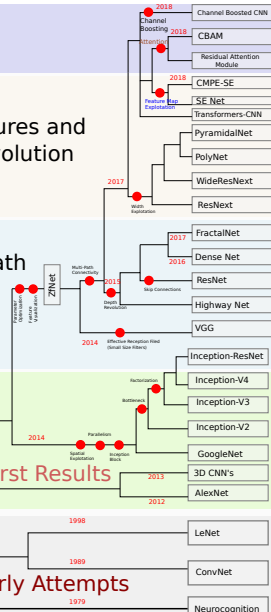
The Revolution



First Results



Early Attempts



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

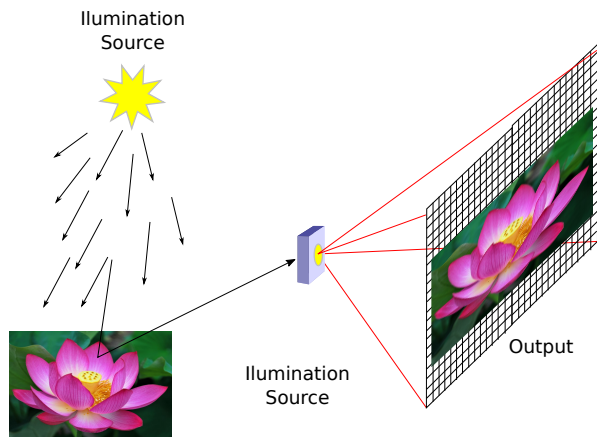
3 Layers

- Convolutional Layer
 - Convolutional Architectures
 - A Little Bit of Notation
 - Deconvolution Layer
 - Alternating Minimization
- 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
- Normalization Layer AKA Batch Normalization
- 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

Something to think about

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

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



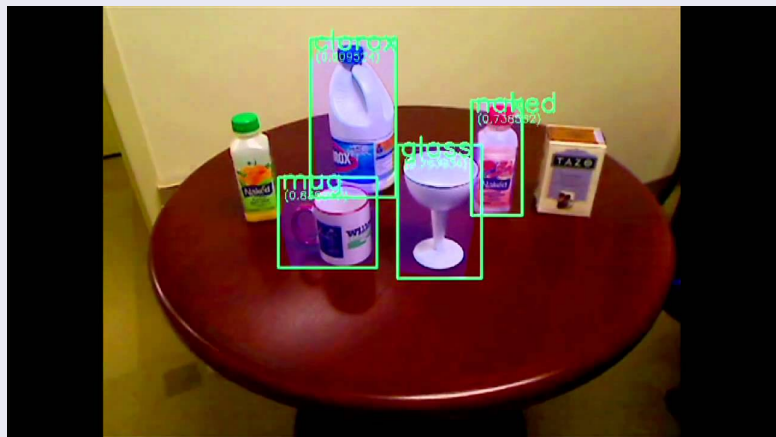
Then

Mid Level Process

Input	Processes	Output
Image	Object Recognition	Attributes
	Segmentation	

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

Why not to use the data sets

- By using a Neural Networks that replicates the 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


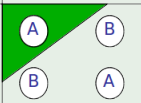


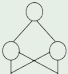
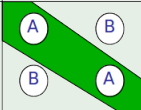
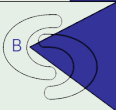

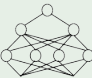
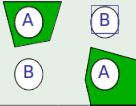

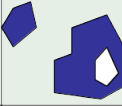
- Convolutional Layer
 - Convolutional Architectures
 - A Little Bit of Notation
 - Deconvolution Layer
 - Alternating Minimization
- Non-Linearity Layer
 - Fixing the Problem, ReLu function
 - Back to the Non-Linearity Layer
- Rectification Layer
- Local Contrast Normalization Layer
- Sub-sampling and Pooling
 - Strides
- Normalization Layer AKA Batch Normalization
- 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

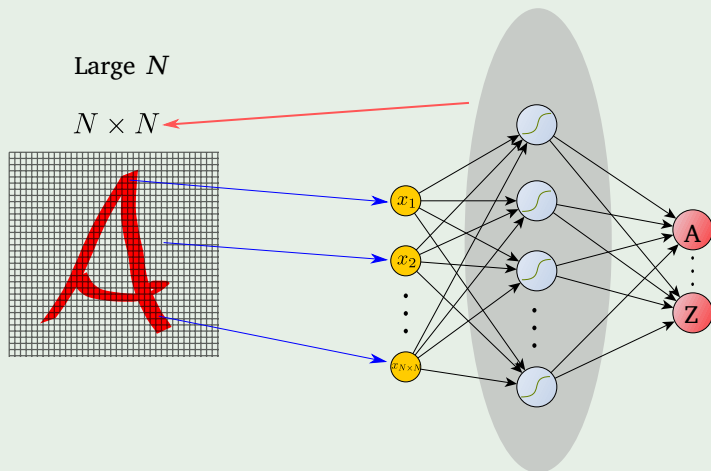
- Convolutional Layer
 - Convolutional Architectures
 - A Little Bit of Notation
 - Deconvolution Layer
 - Alternating Minimization
- Non-Linearity Layer
 - Fixing the Problem, ReLu function
 - Back to the Non-Linearity Layer
- Rectification Layer
- Local Contrast Normalization Layer
- Sub-sampling and Pooling
 - Strides
- Normalization Layer AKA Batch Normalization
- Finally, The Fully Connected Layer

4 An Example of CNN

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

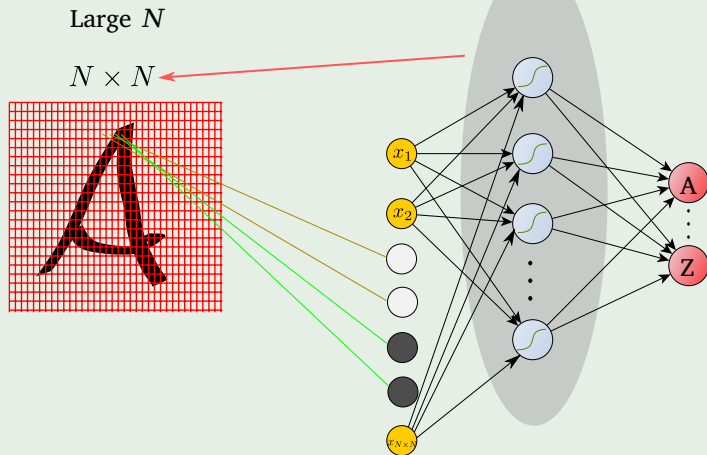
Drawbacks of previous neural networks

The number of trainable parameters becomes extremely large



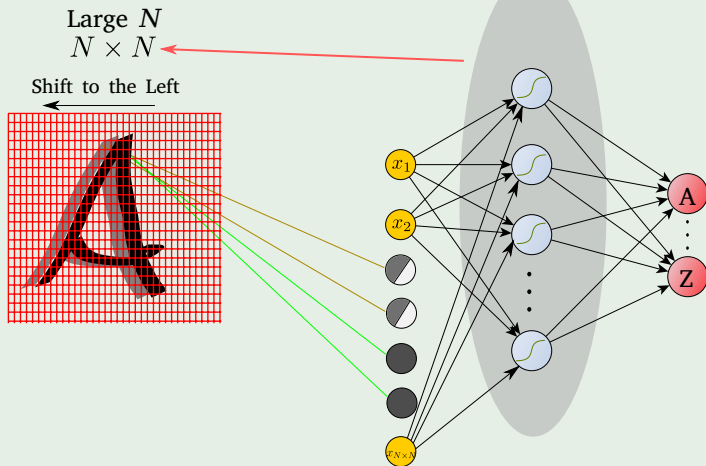
Drawbacks of previous neural networks

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



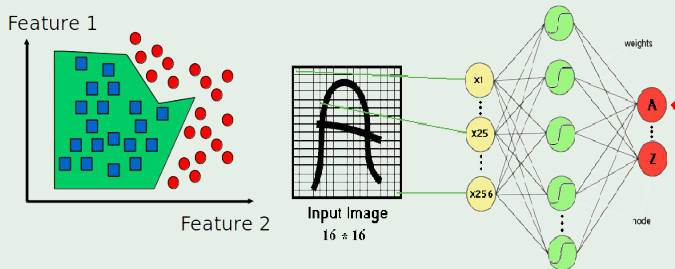
Drawbacks of previous neural networks

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



Drawbacks of previous neural networks

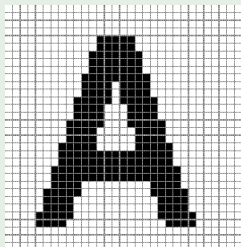
The topology of the input data is completely ignored



For Example

We have

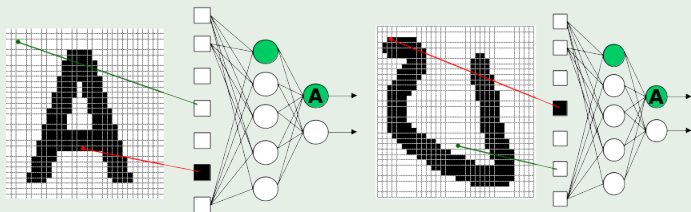
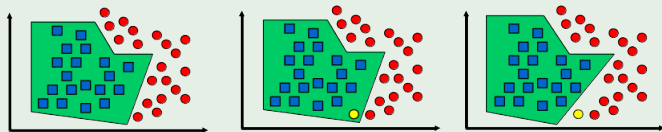
- Black and white patterns: $2^{32 \times 32} = 2^{1024}$
- Gray scale patterns: $256^{32 \times 32} = 256^{1024}$



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.

Problem!!!

- Training time
- Network size
- Free parameters

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
 - Deconvolution Layer
 - Alternating Minimization
- 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
- Normalization Layer AKA Batch Normalization
- 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)

They commented:

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

Hubel/Wiesel Architecture

Something Notable [4]

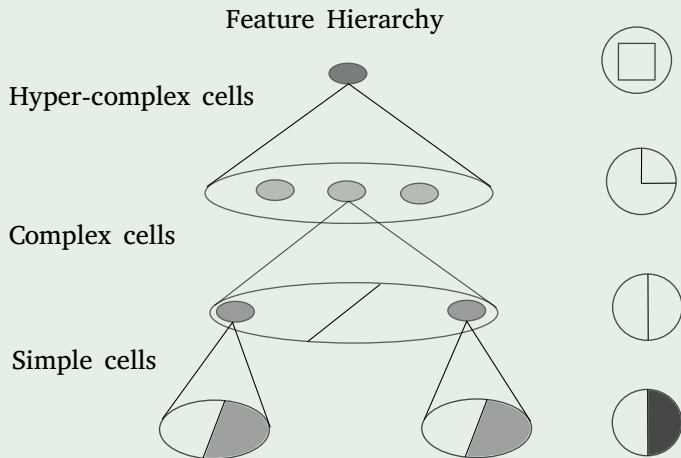
- D. Hubel and T. Wiesel (1959, 1962, Nobel Prize 1981)

They commented

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

Something Like

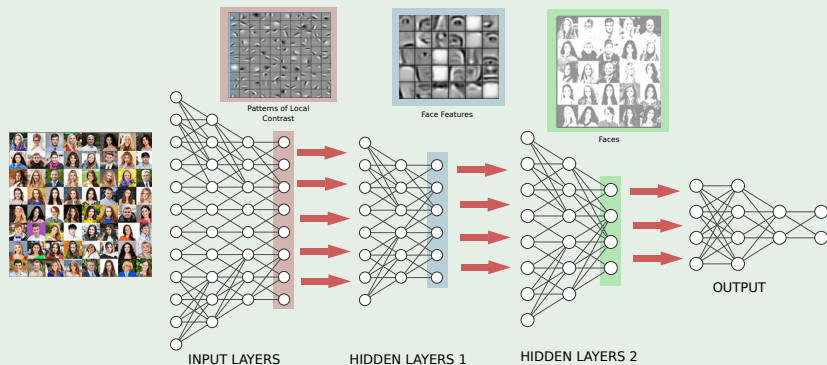
We have



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.

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

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

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.
- 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.

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.

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.

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
 - Deconvolution Layer
 - Alternating Minimization
- 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
- Normalization Layer AKA Batch Normalization
- Finally, The Fully Connected Layer

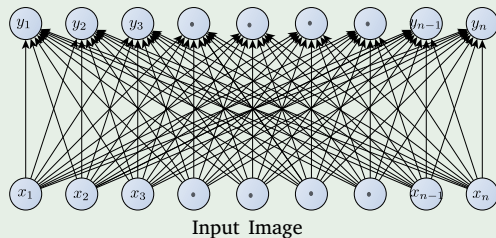
4 An Example of CNN

- 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...



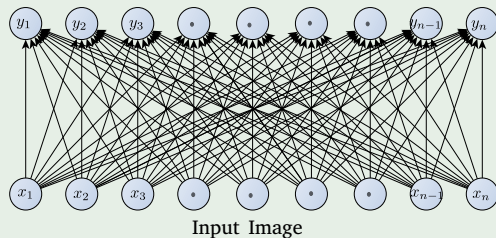
We would have something like this

$$y_i = f \left(\sum_{j=1}^n w_{ij} x_j \right) \quad (1)$$

Local Connectivity

We have the following idea [6]

- Instead of using a full connectivity...



We would have something like this

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

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:

Local Connectivity

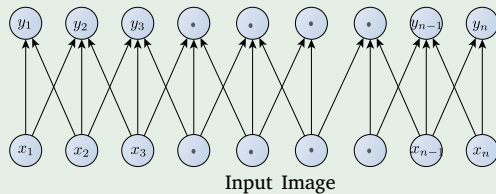
We decide only to connect the neurons in a local way

- 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

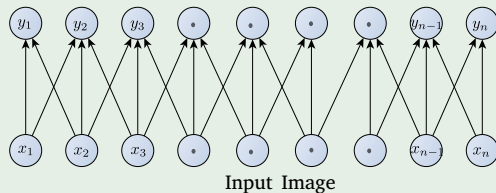


Then, our formula changes

$$y_i = f \left(\sum_{j \in L_p} w_{ij} x_j \right) \quad (2)$$

Example

For gray scale, we get something like this

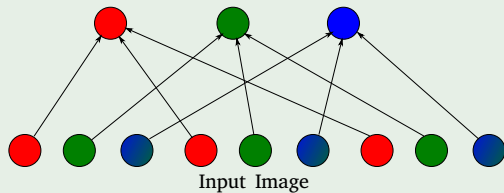


Then, our formula changes

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

Example

In the case of the 3 channels

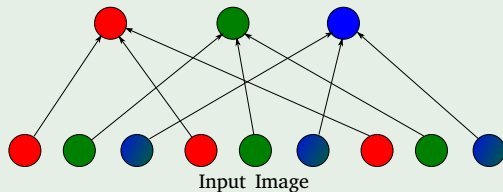


This

$$y_i = f \left(\sum_{i \in L_{p,c}} w_{ij} x_j \right) \quad (3)$$

Example

In the case of the 3 channels



Thus

$$y_i = f \left(\sum_{i \in L_{p,c}} w_i x_i^c \right) \quad (3)$$

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

Solving the following problems...

First

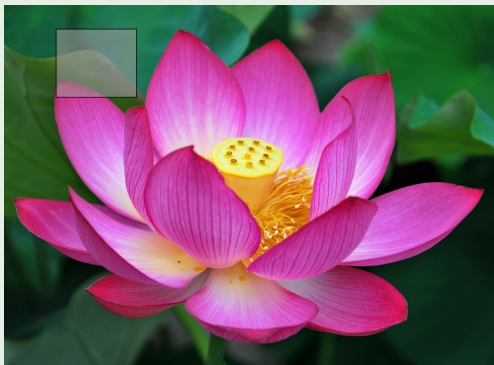
- Fully connected hidden layer would have an unmanageable number of parameters

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
 - Deconvolution Layer
 - Alternating Minimization
- 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
- Normalization Layer AKA Batch Normalization
- 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.

These units are organized into

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

Parameter Sharing

Second Idea

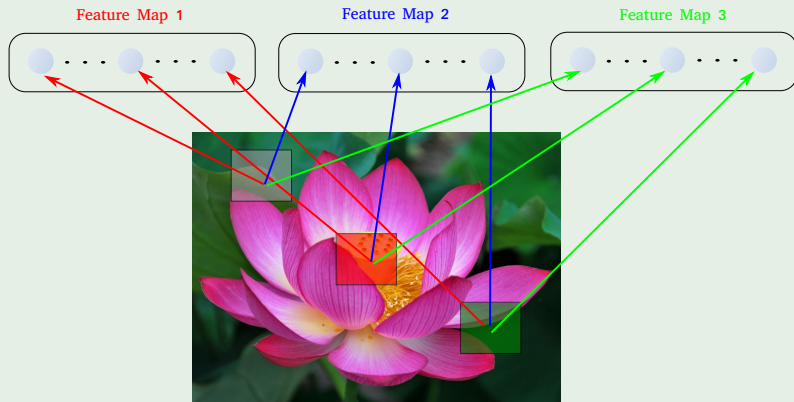
Share matrix of parameters across certain units.

These units are organized into

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

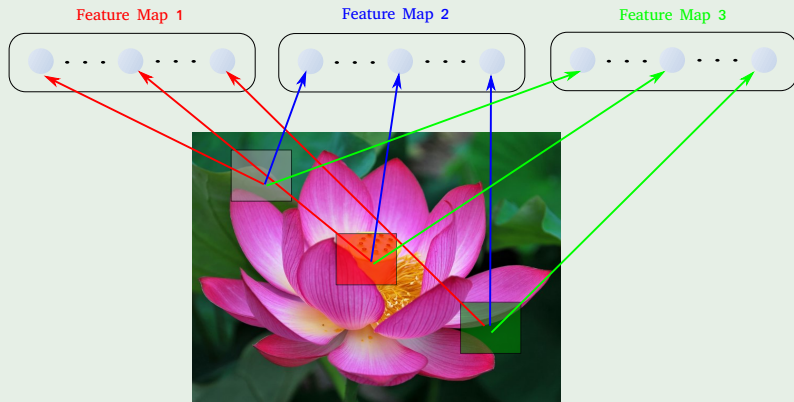
Example

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.
- In each cell of these matrices is the weight to be multiplied with the local input to the local neuron.

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 of 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
 - Deconvolution Layer
 - Alternating Minimization
- 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
- Normalization Layer AKA Batch Normalization
- 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:

- 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
(e.g. $I(x, y) = [0, 255]$)

				$j \rightarrow$								
					79	5	6	90	12	34	2	1
					8	90	12	34	26	78	34	5
		$i \downarrow$			8	1	3	90	12	34	11	61
					77	90	12	34	200	2	9	45
					1	3	90	12	20	1	6	23

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 (with $i \in [0, 255]$):

		$j \rightarrow$							
		79	5	6	90	12	34	2	1
		8	90	12	34	26	78	34	5
	$i \downarrow$	8	1	3	90	12	34	11	61
		77	90	12	34	200	2	9	45
		1	3	90	12	20	1	6	23

Digital Images

In computer vision [2, 7]

We usually operate on digital (discrete) images:

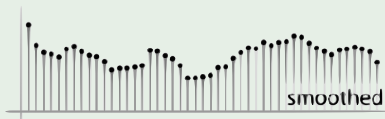
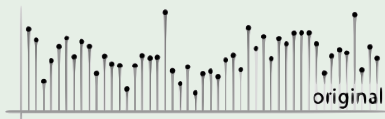
- Sample the 2D space on a regular grid.
- Quantize each sample (round to nearest integer).

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

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

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

$$(I * k)(i) = \sum_{j=-n}^n I(i-j) \times K(j) = \frac{1}{N} \sum_{j=m}^{-m} I(i-j) \quad (4)$$

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	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											

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

$$(I * K)(i, j) = \sum_{s=-m}^{-m} \sum_{l=-m}^m I(i - s, j - l) \times K(s, l) \quad (5)$$

For example, the Box Filter

$$K = \frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} \quad \text{"The Box Filter"} \quad (6)$$

Moving average in 2D

Basically in 2D

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

$$(I * K)(i, j) = \sum_{s=-m}^{-m} \sum_{l=-m}^{m} I(i - s, j - l) \times K(s, l) \quad (5)$$

For example, the Box Filter

$$K = \frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} \quad \text{"The Box Filter"} \quad (6)$$

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

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 Convoluting 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

$K(W, H, C_{in}, C_{out})$

W = Width

H = Height

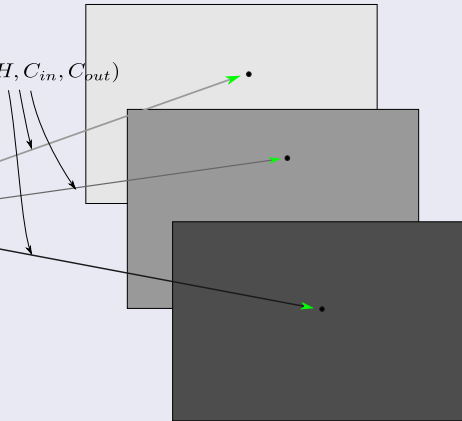
C_{in} = channels input

C_{out} = channel output



Feature Maps

$K(W, H, C_{in}, C_{out})$



Mathematically, we have the following

Map i

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

Therefore

- The convolution works as a
 - ▶ Filter
 - ▶ Encoder
 - ▶ Decoder
 - ▶ etc

Mathematically, we have the following

Map i

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

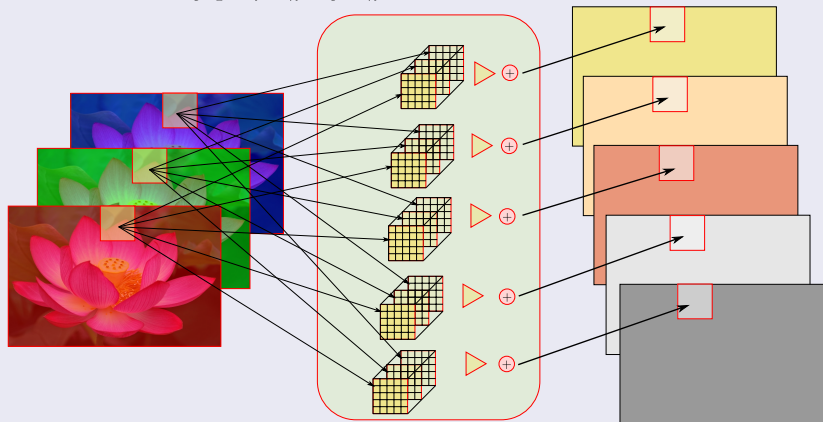
Therefore

- The convolution works as a
 - ▶ 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} .

Therefore

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

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

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

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

Where

- $Y_i^{(l)}$ is the i^{th} feature map in layer l .
- $B_j^{(l)}$ 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, 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 $m_2^{(l-1)} \times m_3^{(l-1)}$.

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

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 $m_2^{(l)} \times m_3^{(l)}$

Something to note

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

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

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

Therefore

The output of layer l

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

Something Notable

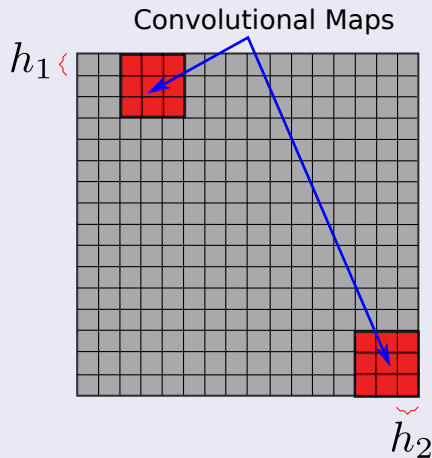
- $m_2^{(l)}$ and $m_3^{(l)}$ are influenced by border effects.
- Therefore, the output feature maps when the Convolutional sum is defined properly have size

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

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

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, the unit at position (x, y) computes

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

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

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

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)

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

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

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

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

They proposed a sparsity idea to start the implementation as

$$C(Y^{(l-1)}) = \sum_{i=1}^{m_1^{(l-1)}} \left\| \sum_{j=1}^{m_1^{(l)}} Y_j^{(l)} * K_{ij}^{(l)} - Y_i^{(l-1)} \right\|_2^2 + \sum_{j=1}^{m_1^{(l)}} |Y_j^{(l)}|^p$$

- Typically, $p = 1$, although other values are possible.

They look for the arguments to minimize a cost of function over a set of images $y = \{y_i\}$

$$\arg \min_{Y_j^{(l)} * K_{ij}^{(l)}} C(y)$$

They proposed a sparsity idea to start the implementation as

$$C(Y^{(l-1)}) = \sum_{i=1}^{m_1^{(l-1)}} \left\| \sum_{j=1}^{m_1^{(l)}} Y_j^{(l)} * K_{ij}^{(l)} - Y_i^{(l-1)} \right\|_2^2 + \sum_{j=1}^{m_1^{(l)}} |Y_j^{(l)}|^p$$

- Typically, $p = 1$, although other values are possible.

They look for the arguments to minimize a cost of function over a set of images $y = \{y^1, \dots, y^I\}$

$$\arg \min_{Y_j^{(l)} * K_{ij}^{(l)}} C(y)$$

Here

Then, we can generalize such cost function for that total set of images (Minbatch)

$$C_l(y) = \frac{\lambda}{2} \sum_{k=1}^I \sum_{i=1}^{m_1^{(l-1)}} \left\| \sum_{j=1}^{m_1^{(l)}} g_{ij}^{(l)} \left(Y_j^{(l,k)} * K_{ij}^{(l)} \right) - Y_i^{(l-1,k)} \right\|_2^2 + \sum_{j=1}^{m_1^{(l)}} \left| Y_j^{(l,k)} \right|^p$$

Here, we have

- $Y_i^{(l-1,k)}$ are the feature maps from the previous layer
- $g_{ij}^{(l)}$ is a fixed binary matrix that determines the connectivity between feature maps at different layers
 - If $Y_j^{(l,k)}$ is connected to certain $Y_i^{(l-1,k)}$ elements

Here

Then, we can generalize such cost function for that total set of images (Minbatch)

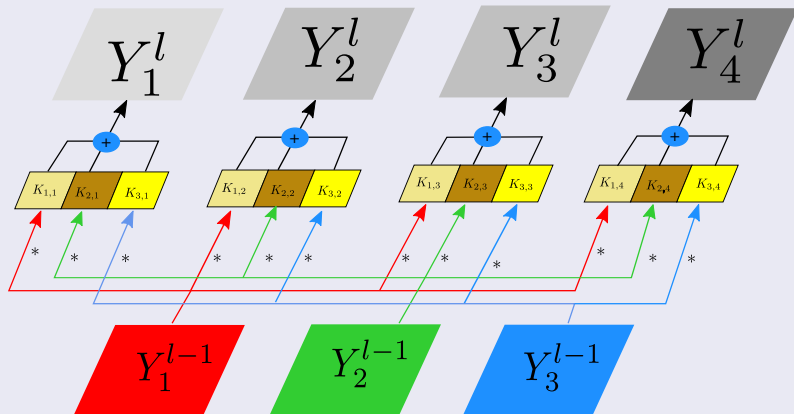
$$C_l(y) = \frac{\lambda}{2} \sum_{k=1}^I \sum_{i=1}^{m_1^{(l-1)}} \left\| \sum_{j=1}^{m_1^{(l)}} g_{ij}^{(l)} \left(Y_j^{(l,k)} * K_{ij}^{(l)} \right) - Y_i^{(l-1,k)} \right\|_2^2 + \sum_{j=1}^{m_1^{(l)}} \left| Y_j^{(l,k)} \right|^p$$

Here, we have

- $Y_i^{(l-1,k)}$ are the feature maps from the previous layer
- $g_{ij}^{(l)}$ is a fixed binary matrix that determines the connectivity between feature maps at different layers
 - ▶ If $Y_j^{(l,k)}$ is connected to certain $Y_i^{(l-1,k)}$ elements

This can be seen as

We have the following layer



They noticed some drawbacks

Using the following optimizations

- Direct Gradient Descent
- Iterative Reweighted Least Squares
- Stochastic Gradient Descent

All of them presented problems!!

- They solved it using a new cost function

They noticed some drawbacks

Using the following optimizations

- Direct Gradient Descent
- Iterative Reweighted Least Squares
- Stochastic Gradient Descent

All of them presented problems!!!

- They solved it using a new cost function

We have that

An interesting use of an auxiliary variable/layer $X_i^{(l,k)}$

$$C_l(y) = \frac{\lambda}{2} \sum_{k=1}^I \sum_{i=1}^{m_1^{(l-1)}} \left\| \sum_{j=1}^{m_1^{(l)}} g_{ij}^{(l)} \left(Y_j^{(l,k)} * K_{ij}^{(l)} \right) - Y_i^{(l-1,k)} \right\|_2^2 + \dots$$
$$\frac{\beta}{2} \sum_{k=1}^I \sum_{j=1}^{m_1^{(l)}} \left\| Y_j^{(l,k)} - X_i^{(l,k)} \right\|_2^2 + \sum_{k=1}^I \sum_{j=1}^{m_1^{(l)}} |Y_j^{(l,k)}|^p$$

This can be solved using

- Alternating minimization...

We have that

An interesting use of an auxiliary variable/layer $X_i^{(l,k)}$

$$C_l(y) = \frac{\lambda}{2} \sum_{k=1}^I \sum_{i=1}^{m_1^{(l-1)}} \left\| \sum_{j=1}^{m_1^{(l)}} g_{ij}^{(l)} \left(Y_j^{(l,k)} * K_{ij}^{(l)} \right) - Y_i^{(l-1,k)} \right\|_2^2 + \dots$$
$$\frac{\beta}{2} \sum_{k=1}^I \sum_{j=1}^{m_1^{(l)}} \left\| Y_j^{(l,k)} - X_i^{(l,k)} \right\|_2^2 + \sum_{k=1}^I \sum_{j=1}^{m_1^{(l)}} |Y_j^{(l,k)}|^p$$

This can be solved using

- Alternating minimization...

This is based on

Fixing the values of $Y_j^{(l,k)}$ and $X_i^{(l,k)}$

- They call these two stages the Y and X sub-problems...

Therefore, they noticed

- These terms introduce the sparsity constraint and gives numerical stability [9, 10]

This is based on

Fixing the values of $Y_j^{(l,k)}$ and $X_i^{(l,k)}$

- They call these two stages the Y and X sub-problems...

Therefore, they noticed

- These terms introduce the sparsity constraint and gives numerical stability [9, 10]

Y sub-problem

Taking the derivative of $Y_j^{(l,k)}$

$$\frac{\partial C_l(y)}{\partial Y_j^{(l,k)}} = \lambda \sum_{i=1}^{m_1^{(l-1)}} F_{ij}^{(l)T} \left[\sum_{t=1}^{m_1^{(l)}} F_{tj}^{(l)} Y_j^{(l,k)} - Y_j^{(l-1,k)} \right] + \beta [Y_j^{(l,k)} - X_j^{(l,k)}] = 0$$

Where

$$F_{ij}^{(l)} = \begin{cases} \text{It is a sparse convolution matrix} & \text{if } g_{ij}^{(l)} = 1 \\ 0 & \text{if } g_{ij}^{(l)} = 0 \end{cases}$$

Y sub-problem

Taking the derivative of $Y_j^{(l,k)}$

$$\frac{\partial C_l(y)}{\partial Y_j^{(l,k)}} = \lambda \sum_{i=1}^{m_1^{(l-1)}} F_{ij}^{(l)T} \left[\sum_{t=1}^{m_1^{(l)}} F_{tj}^{(l)} Y_j^{(l,k)} - Y_j^{(l-1,k)} \right] + \beta [Y_j^{(l,k)} - X_j^{(l,k)}] = 0$$

Where

$$F_{ij}^{(l)} = \begin{cases} \text{It is a sparse convolution matrix} & \text{if } g_{ij}^{(l)} = 1 \\ 0 & \text{if } g_{ij}^{(l)} = 0 \end{cases}$$

Therefore

$F_{ij}^{(l)}$ as a sparse convolution matrix

- Equivalent to convolve with $K_{ij}^{(l)}$

Actually if you fix x , you finish with a linear system $Ax = b$

- Please take a look at the paper... it is interesting
 - ▶ Actually this seems to be the implementation at the Tensorflow framework

Therefore

$F_{ij}^{(l)}$ as a sparse convolution matrix

- Equivalent to convolve with $K_{ij}^{(l)}$

Actually if you fix i , you finish with a linear system $Ax = 0$

- Please take a look at the paper... it is interesting
 - ▶ Actually this seems to be the implementation at the Tensorflow framework

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
 - Deconvolution Layer
 - Alternating Minimization
- **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
- Normalization Layer AKA Batch Normalization
- 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

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

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

$$y(A) = f_l \circ f_{l-1} \circ \dots \circ f_2 \circ f_1(A)$$

With f_l is the last layer.

Therefore, we finish with a sequence of derivatives

$$\frac{\partial y(A)}{\partial w_{li}} = \frac{\partial f_l(f_{l-1})}{\partial f_{l-1}} \cdot \frac{\partial f_{l-1}(f_{l-2})}{\partial f_{l-2}} \cdot \dots \cdot \frac{\partial f_2(f_1)}{\partial f_2} \cdot \frac{\partial f_1(A)}{\partial w_{li}}$$

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

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

With f_t is the last layer.

Therefore we implement a sequence of derivative:

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

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

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

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, 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

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

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

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, 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

Therefore

The maximum for the derivative of the sigmoid

- $f'(0) = 0.25$

Therefore, Given a **Deep** Convolutional Network

- We could finish with

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

Vanishing derivative

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

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

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

It is called ReLU or Rectifier

With a smooth approximation (Softplus function)

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

Thus

The need to introduce a new function

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

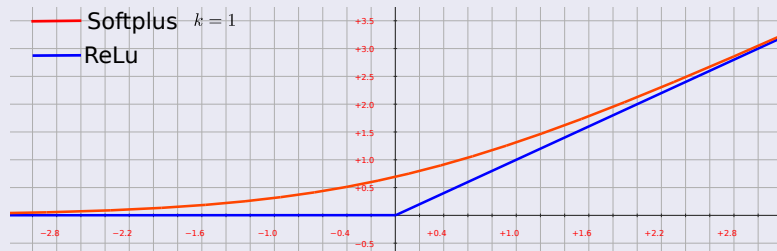
It is called ReLu or Rectifier

With a smooth approximation (Softplus function)

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

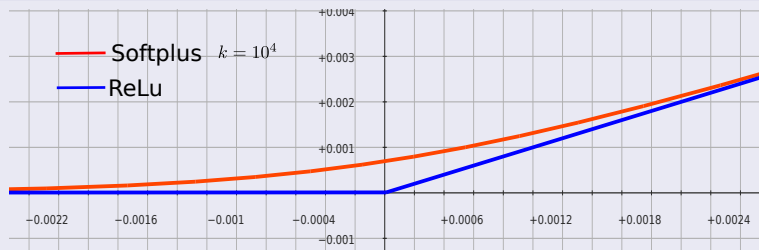
Therefore, we have

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.

What about the output

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

Each of them of size

$$m_2^{(l-1)} \times m_3^{(l-1)} \quad (7)$$

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

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

$$m_2^{(l-1)} \times m_3^{(l-1)} \quad (7)$$

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

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

$$m_2^{(l-1)} \times m_3^{(l-1)} \quad (7)$$

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

Thus

With the final output

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

Where

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

You can also add a gain to compensate

$$Y_i^{(l)} = g_i f \left(Y_i^{(l-1)} \right) \quad (9)$$

Thus

With the final output

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

Where

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

You can also add a gain to compensate

$$Y_i^{(l)} = g_i f \left(Y_i^{(l-1)} \right) \quad (9)$$

Thus

With the final output

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

Where

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

You can also add a gain to compensate

$$Y_i^{(l)} = g_i f \left(Y_i^{(l-1)} \right) \quad (9)$$

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
 - Deconvolution Layer
 - Alternating Minimization
- 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
- Normalization Layer AKA Batch Normalization
- 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 $m_2^{(l-1)} \times m_3^{(l-1)}$.

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

$$Y_i^{(l)} = |Y_i^{(l-1)}| \quad (10)$$

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.

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

$$Y_i^{(l)} = |Y_i^{(l)}| \quad (10)$$

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.

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

$$Y_i^{(l)} = |Y_i^{(l)}| \quad (10)$$

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.

You can find this in

K. Jarrett, K. Kavukcuoglu, 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.

Thus

We have that

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

You can find this in

K. Jarrett, K. Kavukcuoglu, 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.

Thus

We have that

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

You can find this in

K. Jarrett, K. Kavukcuoglu, 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?

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

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

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 other way

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

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
 - Deconvolution Layer
 - Alternating Minimization
- 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
- Normalization Layer AKA Batch Normalization
- 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:

- 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.
- Brightness 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.

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.

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.

With the operation

$$Y_i^{(l)} = Y_i^{(l-1)} - \sum_{j=1}^{m_1^{(l-1)}} K_{G(\sigma)} * Y_j^{(l-1)} \quad (11)$$

With

$$\left(K_{G(\sigma)}\right)_{x,y} = \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left\{-\frac{x^2 + y^2}{2\sigma^2}\right\} \quad (12)$$

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

$$Y_i^{(l)} = Y_i^{(l-1)} - \sum_{j=1}^{m_1^{(l-1)}} K_{G(\sigma)} * Y_j^{(l-1)} \quad (11)$$

$$\left(K_{G(\sigma)}\right)_{x,y} = \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left\{-\frac{x^2 + y^2}{2\sigma^2}\right\} \quad (12)$$

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

$$Y_i^{(l)} = Y_i^{(l-1)} - \sum_{j=1}^{m_1^{(l-1)}} K_{G(\sigma)} * Y_j^{(l-1)} \quad (11)$$

With

$$\left(K_{G(\sigma)}\right)_{x,y} = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{x^2 + y^2}{2\sigma^2}\right\} \quad (12)$$

Brightness Normalization

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

$$\left(Y_i^{(l)}\right)_{x,y} = \frac{\left(Y_i^{(l-1)}\right)_{x,y}}{\left(\kappa + \lambda \sum_{j=1}^{m_1^{(l-1)}} \left(Y_j^{(l-1)}\right)_{x,y}^2\right)^\mu} \quad (13)$$

Where

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

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

validation set.

Brightness Normalization

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

$$\left(Y_i^{(l)}\right)_{x,y} = \frac{\left(Y_i^{(l-1)}\right)_{x,y}}{\left(\kappa + \lambda \sum_{j=1}^{m_1^{(l-1)}} \left(Y_j^{(l-1)}\right)_{x,y}^2\right)^\mu} \quad (13)$$

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
 - Deconvolution Layer
 - Alternating Minimization
- 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
- Normalization Layer AKA Batch Normalization
- 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.

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 Layer

Motivation

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

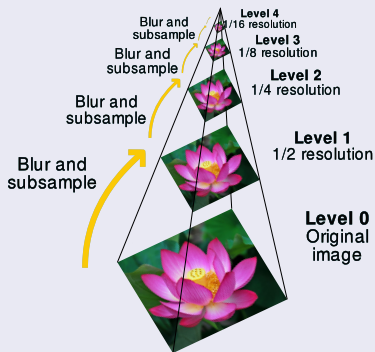
How?

- 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 know sub-sampling pyramid



How is sub-sampling implemented?

We know that Image Pyramids

They were designed to find:

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

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

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

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

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

The basic idea is to skip a fixed number of pixels

Therefore the size of the output feature map is given by

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

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 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.

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.

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 $(l - 1)$.

Thus

In the previous example

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

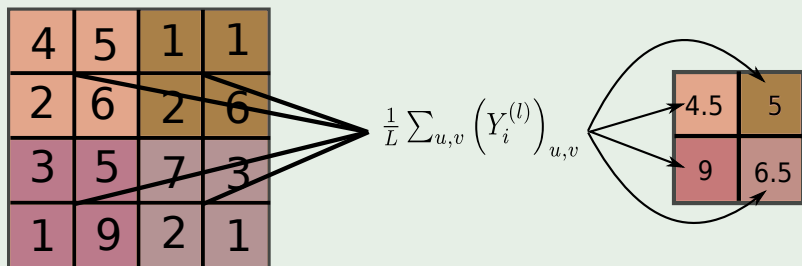
Each unit

- 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

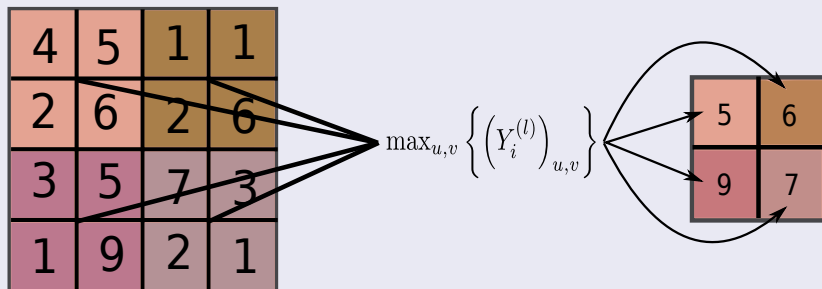
When using a boxcar filter, the operation is called average pooling and the layer denoted by P_A .



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 [11].

The small amount

- In the case of the previous examples, 1 pixel

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 [11].

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

However, we have another way of doing pooling

- Striding!!!

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

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

There is a problem...

- 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...

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

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

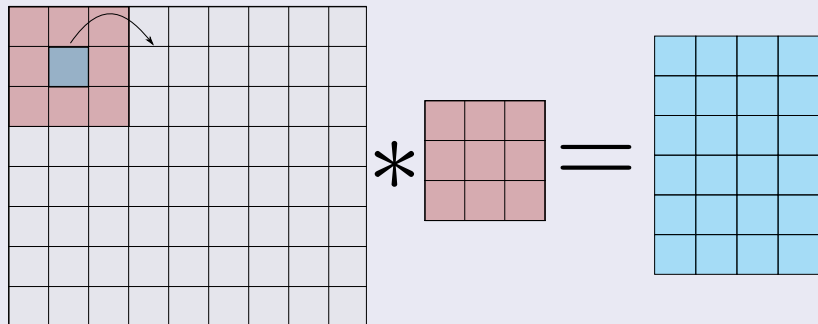
This is a Heuristic ...

- 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

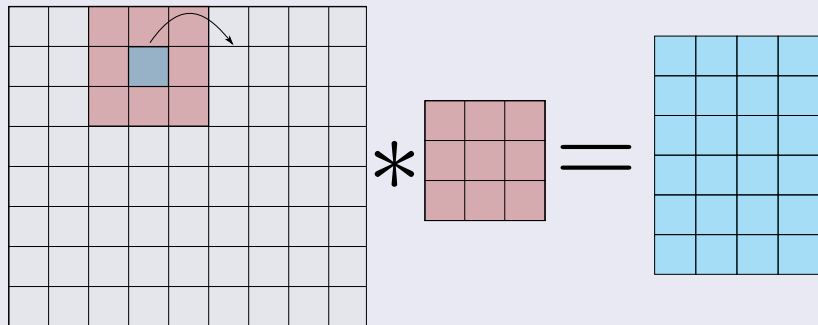
Horizontal Stride $r = 2$



Example

Horizontal Stride

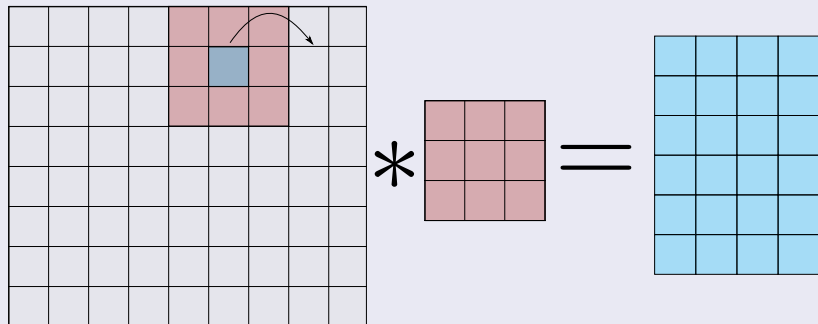
Horizontal Stride $r = 2$



Example

Horizontal Stride

Horizontal Stride $r = 2$



There are attempts to understand its effects

At Convolution Level and using Tensors [13]

- “Take it in your stride: Do we need striding in CNNs?” by Chen Kong, Simon Lucey [14]

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

- You need a little bit of notation...

There are attempts to understand its effects

At Convolution Level and using Tensors [13]

- “Take it in your stride: Do we need striding in CNNs?” by Chen Kong, Simon Lucey [14]

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
 - Deconvolution Layer
 - Alternating Minimization
- 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
- **Normalization Layer AKA Batch Normalization**
- Finally, The Fully Connected Layer

4 An Example of CNN

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

Here, the people at Google [15] around 2015

They commented in the “Internal Covariate Shift Phenomena”

- Due to the change in the distribution of each layer’s input

What claim

- The min-batch forces to have those changes which impact on the learning capabilities of the network.

In Neural Networks, they define this

- Internal Covariate Shift as the change in the distribution of network activation’s due to the change in network parameters during training.

Here, the people at Google [15] around 2015

They commented in the “Internal Covariate Shift Phenomena”

- Due to the change in the distribution of each layer's input

They claim

- The min-batch forces to have those changes which impact on the learning capabilities of the network.

In Neural Networks they define this

- Internal Covariate Shift as the change in the distribution of network activation's due to the change in network parameters during training.

Here, the people at Google [15] around 2015

They commented in the “Internal Covariate Shift Phenomena”

- Due to the change in the distribution of each layer’s input

They claim

- The min-batch forces to have those changes which impact on the learning capabilities of the network.

In Neural Networks, they define this

- Internal Covariate Shift as the change in the distribution of network activation’s due to the change in network parameters during training.

They gave the following reasons

Consider a layer with the input u that adds the learned bias b

- Then, it normalizes the result by subtracting the mean of the activation over the training data:

$$\hat{\mathbf{x}} = \mathbf{x} - E[\mathbf{x}]$$

- ▶ $\mathcal{X} = \{\mathbf{x}, \dots, \mathbf{x}_N\}$ the data samples and $E[\mathbf{x}] = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$

Now, if the gradient ignores the dependence of $E[\mathbf{x}]$ on b

- Then $b = b + \Delta b$ where $\Delta b \propto -\frac{\partial l}{\partial x}$

Finally

$$u + (b + \Delta b) - E[u + (b + \Delta b)] = u + b - E[u + b]$$

They gave the following reasons

Consider a layer with the input u that adds the learned bias b

- Then, it normalizes the result by subtracting the mean of the activation over the training data:

$$\hat{\mathbf{x}} = \mathbf{x} - E[\mathbf{x}]$$

- ▶ $\mathcal{X} = \{\mathbf{x}, \dots, \mathbf{x}_N\}$ the data samples and $E[\mathbf{x}] = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$

Now, if the gradient ignores the dependence of $E[\mathbf{x}]$ on b

- Then $b = b + \Delta b$ where $\Delta b \propto -\frac{\partial l}{\partial x}$

$$u + (b + \Delta b) - E[u + (b + \Delta b)] = u + b - E[u + b]$$

They gave the following reasons

Consider a layer with the input u that adds the learned bias b

- Then, it normalizes the result by subtracting the mean of the activation over the training data:

$$\hat{\mathbf{x}} = \mathbf{x} - E[\mathbf{x}]$$

- ▶ $\mathcal{X} = \{\mathbf{x}, \dots, \mathbf{x}_N\}$ the data samples and $E[\mathbf{x}] = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$

Now, if the gradient ignores the dependence of $E[\mathbf{x}]$ on b

- Then $b = b + \Delta b$ where $\Delta b \propto -\frac{\partial l}{\partial x}$

Finally

$$u + (b + \Delta b) - E[u + (b + \Delta b)] = u + b - E[u + b]$$

Then

The following will happen

- The update to b by Δb leads to **no change** in the output of the layer.

Therefore

- We need to integrate the normalization into the process of training.

Then

The following will happen

- The update to b by Δb leads to **no change** in the output of the layer.

Therefore

- We need to integrate the normalization into the process of training.

Normalization via Mini-Batch Statistic

It is possible to describe the **normalization** as a transformation layer

$$\hat{\mathbf{x}} = \text{Norm}(\mathbf{x}, \mathcal{X})$$

- Which depends on all the training samples \mathcal{X} which also depends on the layer parameters

For back-propagation, we will need to generate the following terms

$$\frac{\partial \text{Norm}(\mathbf{x}, \mathcal{X})}{\partial \mathbf{x}} \quad \text{and} \quad \frac{\partial \text{Norm}(\mathbf{x}, \mathcal{X})}{\partial \mathcal{X}}$$

Normalization via Mini-Batch Statistic

It is possible to describe the **normalization** as a transformation layer

$$\hat{\mathbf{x}} = \text{Norm}(\mathbf{x}, \mathcal{X})$$

- Which depends on all the training samples \mathcal{X} which also depends on the layer parameters

For back-propagation, we will need to generate the following terms

$$\frac{\partial \text{Norm}(\mathbf{x}, \mathcal{X})}{\partial \mathbf{x}} \quad \text{and} \quad \frac{\partial \text{Norm}(\mathbf{x}, \mathcal{X})}{\partial \mathcal{X}}$$

Definition of Whitening

Whitening

- Suppose X is a random (column) vector with non-singular covariance matrix Σ and mean 0.

Then

- Then the transformation $Y = WX$ with a whitening matrix W satisfying the condition $W^T W = \Sigma^{-1}$ yields the whitened random vector Y with unit diagonal covariance.

Definition of Whitening

Whitening

- Suppose X is a random (column) vector with non-singular covariance matrix Σ and mean 0.

Then

- Then the transformation $Y = WX$ with a whitening matrix W satisfying the condition $W^T W = \Sigma^{-1}$ yields the whitened random vector Y with unit diagonal covariance.

Such Normalization

It could be used for all layer

- But whitening the layer inputs is expensive, as it requires computing the covariance matrix

$$\text{Cov}[\mathbf{x}] = E_{\mathbf{x} \in \mathcal{X}} [\mathbf{x}\mathbf{x}^T] \quad \text{and} \quad E[\mathbf{x}] E[\mathbf{x}]^T$$

- ▶ To produce the whitened activations

Therefore

A Better Options, we can normalize each input layer

$$\hat{\mathbf{x}}^{(k)} = \frac{\mathbf{x}^{(k)} - \mu}{\sigma}$$

- with $\mu = E[\mathbf{x}^{(k)}]$ and $\sigma^2 = Var[\mathbf{x}^{(k)}]$

This allows to speed up convergence

- Simply normalizing each input of a layer may change what the layer can represent.

So, we need to insert a transformation in the network

- Which can represent the identity transform

Therefore

A Better Options, we can normalize each input layer

$$\hat{\mathbf{x}}^{(k)} = \frac{\mathbf{x}^{(k)} - \mu}{\sigma}$$

- with $\mu = E[\mathbf{x}^{(k)}]$ and $\sigma^2 = Var[\mathbf{x}^{(k)}]$

This allows to speed up convergence

- Simply normalizing each input of a layer may change what the layer can represent.

So, we need to insert a transformation in the network

- Which can represent the identity transform

Therefore

A Better Options, we can normalize each input layer

$$\hat{\mathbf{x}}^{(k)} = \frac{\mathbf{x}^{(k)} - \mu}{\sigma}$$

- with $\mu = E[\mathbf{x}^{(k)}]$ and $\sigma^2 = Var[\mathbf{x}^{(k)}]$

This allows to speed up convergence

- Simply normalizing each input of a layer may change what the layer can represent.

So, we need to insert a transformation in the network

- Which can represent the identity transform

The Transformation

The Linear transformation

$$\mathbf{y}^{(k)} = \gamma^{(k)} \hat{\mathbf{x}}^{(k)} + \beta^{(k)}$$

The parameters

- This allow to recover the identity by setting $\gamma^{(k)} = \sqrt{Var [\mathbf{x}^{(k)}]}$ and $\beta^{(k)} = E [\mathbf{x}^{(k)}]$ if necessary.

The Transformation

The Linear transformation

$$\mathbf{y}^{(k)} = \gamma^{(k)} \hat{\mathbf{x}}^{(k)} + \beta^{(k)}$$

The parameters $\gamma^{(k)}, \beta^{(k)}$

- This allow to recover the identity by setting $\gamma^{(k)} = \sqrt{Var [\mathbf{x}^{(k)}]}$ and $\beta^{(k)} = E [\mathbf{x}^{(k)}]$ if necessary.

Finally

Batch Normalizing Transform

Input: Values of \mathbf{x} over a mini-batch: $\mathcal{B} = \{\mathbf{x}_{1\dots m}\}$, Parameters to be learned: γ, β

Output: $\{y_i = BN_{\gamma, \beta}(\mathbf{x}_i)\}$

$$\textcircled{1} \mu_{\mathcal{B}} = \frac{1}{m} \sum_{i=1}^m x_i$$

$$\textcircled{2} \sigma_{\mathcal{B}}^2 = \frac{1}{m} \sum_{i=1}^m (x_i - \mu)^2$$

$$\textcircled{3} \hat{x} = \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$$

$$\textcircled{4} y_i = \gamma^{(k)} \hat{x}_i + \beta = BN_{\gamma, \beta}(\mathbf{x}_i)$$

Finally

Batch Normalizing Transform

Input: Values of \mathbf{x} over a mini-batch: $\mathcal{B} = \{\mathbf{x}_{1\dots m}\}$, Parameters to be learned: γ, β

Output: $\{y_i = BN_{\gamma, \beta}(\mathbf{x}_i)\}$

① $\mu_{\mathcal{B}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i$

② $\sigma_{\mathcal{B}}^2 = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}_i - \mu)^2$

③ $\hat{\mathbf{x}} = \frac{\mathbf{x}_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$

④ $y_i = \gamma^{(k)} \hat{\mathbf{x}}_i + \beta = BN_{\gamma, \beta}(\mathbf{x}_i)$

Finally

Batch Normalizing Transform

Input: Values of \mathbf{x} over a mini-batch: $\mathcal{B} = \{\mathbf{x}_{1\dots m}\}$, Parameters to be learned: γ, β

Output: $\{y_i = BN_{\gamma, \beta}(\mathbf{x}_i)\}$

1 $\mu_{\mathcal{B}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i$

2 $\sigma_{\mathcal{B}}^2 = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}_i - \mu)^2$

3 $\hat{\mathbf{x}} = \frac{\mathbf{x}_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$

4 $y_i = \gamma^{(k)} \hat{\mathbf{x}}_i + \beta = BN_{\gamma, \beta}(\mathbf{x}_i)$

Finally

Batch Normalizing Transform

Input: Values of \mathbf{x} over a mini-batch: $\mathcal{B} = \{\mathbf{x}_{1\dots m}\}$, Parameters to be learned: γ, β

Output: $\{y_i = BN_{\gamma, \beta}(\mathbf{x}_i)\}$

① $\mu_{\mathcal{B}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i$

② $\sigma_{\mathcal{B}}^2 = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}_i - \mu)^2$

③ $\hat{\mathbf{x}} = \frac{\mathbf{x}_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$

④ $y_i = \gamma^{(k)} \hat{\mathbf{x}}_i + \beta = BN_{\gamma, \beta}(\mathbf{x}_i)$

Finally

Batch Normalizing Transform

Input: Values of \mathbf{x} over a mini-batch: $\mathcal{B} = \{\mathbf{x}_{1\dots m}\}$, Parameters to be learned: γ, β

Output: $\{y_i = BN_{\gamma, \beta}(\mathbf{x}_i)\}$

- 1 $\mu_{\mathcal{B}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i$
- 2 $\sigma_{\mathcal{B}}^2 = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}_i - \mu)^2$
- 3 $\hat{\mathbf{x}} = \frac{\mathbf{x}_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$
- 4 $\mathbf{y}_i = \gamma^{(k)} \hat{\mathbf{x}}_i + \beta = BN_{\gamma, \beta}(\mathbf{x}_i)$

Backpropagation

We have the following equations by using the loss function l

$$\textcircled{1} \quad \frac{\partial l}{\partial \hat{x}_i} = \frac{\partial l}{\partial y_i} \times \gamma$$

$$\textcircled{2} \quad \frac{\partial l}{\partial \sigma_B^2} = \sum_{i=1}^m \frac{\partial l}{\partial x_i} \times (x_i - \mu_B) \times \left(-\frac{1}{2}\right) \times (\sigma_B^2 + \epsilon)^{-\frac{3}{2}}$$

$$\textcircled{3} \quad \frac{\partial l}{\partial \mu_B} = \left(\sum_{i=1}^m \frac{\partial l}{\partial x_i} \times \frac{-1}{\sqrt{\sigma_B^2 + \epsilon}} \right) + \frac{\partial l}{\partial \sigma_B^2} \times \frac{\sum_{i=1}^m -2 \times (x_i - \mu_B)}{m}$$

$$\textcircled{4} \quad \frac{\partial l}{\partial x_i} = \frac{\partial l}{\partial x_i} \times \frac{1}{\sqrt{\sigma_B^2 + \epsilon}} + \frac{\partial l}{\partial \sigma_B^2} \times \frac{2 \times (x_i - \mu_B)}{m} + \frac{\partial l}{\partial \mu_B} \times \frac{1}{m}$$

$$\textcircled{5} \quad \frac{\partial l}{\partial \gamma} = \sum_{i=1}^m \frac{\partial l}{\partial y_i} \times \hat{x}_i$$

$$\textcircled{6} \quad \frac{\partial l}{\partial \beta} = \sum_{i=1}^m \frac{\partial l}{\partial y_i}$$

Backpropagation

We have the following equations by using the loss function l

$$\textcircled{1} \quad \frac{\partial l}{\partial \hat{\mathbf{x}}_i} = \frac{\partial l}{\partial \mathbf{y}_i} \times \gamma$$

$$\textcircled{2} \quad \frac{\partial l}{\partial \sigma_B^2} = \sum_{i=1}^m \frac{\partial l}{\partial \hat{\mathbf{x}}_i} \times (\mathbf{x}_i - \mu_B) \times \left(-\frac{1}{2}\right) \times (\sigma_B^2 + \epsilon)^{-\frac{3}{2}}$$

$$\textcircled{3} \quad \frac{\partial l}{\partial \mu_B} = \left(\sum_{i=1}^m \frac{\partial l}{\partial \hat{\mathbf{x}}_i} \times \frac{-1}{\sqrt{\sigma_B^2 + \epsilon}} \right) + \frac{\partial l}{\partial \sigma_B^2} \times \frac{\sum_{i=1}^m -2 \times (\mathbf{x}_i - \mu_B)}{m}$$

$$\textcircled{4} \quad \frac{\partial l}{\partial \mathbf{x}_i} = \frac{\partial l}{\partial \hat{\mathbf{x}}_i} \times \frac{1}{\sqrt{\sigma_B^2 + \epsilon}} + \frac{\partial l}{\partial \sigma_B^2} \times \frac{2 \times (\mathbf{x}_i - \mu_B)}{m} + \frac{\partial l}{\partial \mu_B} \times \frac{1}{m}$$

$$\textcircled{5} \quad \frac{\partial l}{\partial \gamma} = \sum_{i=1}^m \frac{\partial l}{\partial \hat{\mathbf{y}}_i} \times \hat{\mathbf{x}}_i$$

$$\textcircled{6} \quad \frac{\partial l}{\partial \beta} = \sum_{i=1}^m \frac{\partial l}{\partial \hat{\mathbf{y}}_i}$$

Backpropagation

We have the following equations by using the loss function l

$$\textcircled{1} \quad \frac{\partial l}{\partial \hat{\mathbf{x}}_i} = \frac{\partial l}{\partial \mathbf{y}_i} \times \gamma$$

$$\textcircled{2} \quad \frac{\partial l}{\partial \sigma_B^2} = \sum_{i=1}^m \frac{\partial l}{\partial \hat{\mathbf{x}}_i} \times (\mathbf{x}_i - \mu_B) \times \left(-\frac{1}{2}\right) \times (\sigma_B^2 + \epsilon)^{-\frac{3}{2}}$$

$$\textcircled{3} \quad \frac{\partial l}{\partial \mu_B} = \left(\sum_{i=1}^m \frac{\partial l}{\partial \hat{\mathbf{x}}_i} \times \frac{-1}{\sqrt{\sigma_B^2 + \epsilon}} \right) + \frac{\partial l}{\partial \sigma_B^2} \times \frac{\sum_{i=1}^m -2 \times (\mathbf{x}_i - \mu_B)}{m}$$

$$\textcircled{4} \quad \frac{\partial l}{\partial \mathbf{x}_i} = \frac{\partial l}{\partial \hat{\mathbf{x}}_i} \times \frac{1}{\sqrt{\sigma_B^2 + \epsilon}} + \frac{\partial l}{\partial \sigma_B^2} \times \frac{2 \times (\mathbf{x}_i - \mu_B)}{m} + \frac{\partial l}{\partial \mu_B} \times \frac{1}{m}$$

$$\textcircled{5} \quad \frac{\partial l}{\partial \mathbf{y}} = \sum_{i=1}^m \frac{\partial l}{\partial \hat{\mathbf{x}}_i} \times \hat{\mathbf{x}}_i$$

$$\textcircled{6} \quad \frac{\partial l}{\partial \mathbf{p}} = \sum_{i=1}^m \frac{\partial l}{\partial \hat{\mathbf{x}}_i}$$

Backpropagation

We have the following equations by using the loss function l

$$① \quad \frac{\partial l}{\partial \hat{\mathbf{x}}_i} = \frac{\partial l}{\partial \mathbf{y}_i} \times \gamma$$

$$② \quad \frac{\partial l}{\partial \sigma_B^2} = \sum_{i=1}^m \frac{\partial l}{\partial \hat{\mathbf{x}}_i} \times (\mathbf{x}_i - \mu_B) \times \left(-\frac{1}{2}\right) \times (\sigma_B^2 + \epsilon)^{-\frac{3}{2}}$$

$$③ \quad \frac{\partial l}{\partial \mu_B} = \left(\sum_{i=1}^m \frac{\partial l}{\partial \hat{\mathbf{x}}_i} \times \frac{-1}{\sqrt{\sigma_B^2 + \epsilon}} \right) + \frac{\partial l}{\partial \sigma_B^2} \times \frac{\sum_{i=1}^m -2 \times (\mathbf{x}_i - \mu_B)}{m}$$

$$④ \quad \frac{\partial l}{\partial \mathbf{x}_i} = \frac{\partial \hat{\mathbf{x}}_i}{\partial \mathbf{x}_i} \times \frac{1}{\sqrt{\sigma_B^2 + \epsilon}} + \frac{\partial l}{\partial \sigma_B^2} \times \frac{2 \times (\mathbf{x}_i - \mu_B)}{m} + \frac{\partial l}{\partial \mu_B} \times \frac{1}{m}$$

$$\frac{\partial l}{\partial \mathbf{y}_i} = \sum_{i=1}^m \frac{\partial l}{\partial \hat{\mathbf{x}}_i} \times \hat{\mathbf{x}}_i$$

$$\frac{\partial l}{\partial \mu} = \sum_{i=1}^m \frac{\partial l}{\partial \hat{\mathbf{x}}_i}$$

Backpropagation

We have the following equations by using the loss function l

$$① \quad \frac{\partial l}{\partial \hat{\mathbf{x}}_i} = \frac{\partial l}{\partial \mathbf{y}_i} \times \gamma$$

$$② \quad \frac{\partial l}{\partial \sigma_B^2} = \sum_{i=1}^m \frac{\partial l}{\partial \hat{\mathbf{x}}_i} \times (\mathbf{x}_i - \mu_B) \times \left(-\frac{1}{2}\right) \times (\sigma_B^2 + \epsilon)^{-\frac{3}{2}}$$

$$③ \quad \frac{\partial l}{\partial \mu_B} = \left(\sum_{i=1}^m \frac{\partial l}{\partial \hat{\mathbf{x}}_i} \times \frac{-1}{\sqrt{\sigma_B^2 + \epsilon}} \right) + \frac{\partial l}{\partial \sigma_B^2} \times \frac{\sum_{i=1}^m -2 \times (\mathbf{x}_i - \mu_B)}{m}$$

$$④ \quad \frac{\partial l}{\partial \mathbf{x}_i} = \frac{\partial \hat{\mathbf{x}}_i}{\partial \mathbf{x}_i} \times \frac{1}{\sqrt{\sigma_B^2 + \epsilon}} + \frac{\partial l}{\partial \sigma_B^2} \times \frac{2 \times (\mathbf{x}_i - \mu_B)}{m} + \frac{\partial l}{\partial \mu_B} \times \frac{1}{m}$$

$$⑤ \quad \frac{\partial l}{\partial \gamma} = \sum_{i=1}^m \frac{\partial l}{\partial \mathbf{y}_i} \times \hat{\mathbf{x}}_i$$

$$⑥ \quad \frac{\partial l}{\partial \beta} = \sum_{i=1}^m \frac{\partial l}{\partial \mathbf{y}_i}$$

Backpropagation

We have the following equations by using the loss function l

$$① \quad \frac{\partial l}{\partial \hat{\mathbf{x}}_i} = \frac{\partial l}{\partial \mathbf{y}_i} \times \gamma$$

$$② \quad \frac{\partial l}{\partial \sigma_B^2} = \sum_{i=1}^m \frac{\partial l}{\partial \hat{\mathbf{x}}_i} \times (\mathbf{x}_i - \mu_B) \times \left(-\frac{1}{2}\right) \times (\sigma_B^2 + \epsilon)^{-\frac{3}{2}}$$

$$③ \quad \frac{\partial l}{\partial \mu_B} = \left(\sum_{i=1}^m \frac{\partial l}{\partial \hat{\mathbf{x}}_i} \times \frac{-1}{\sqrt{\sigma_B^2 + \epsilon}} \right) + \frac{\partial l}{\partial \sigma_B^2} \times \frac{\sum_{i=1}^m -2 \times (\mathbf{x}_i - \mu_B)}{m}$$

$$④ \quad \frac{\partial l}{\partial \mathbf{x}_i} = \frac{\partial \hat{\mathbf{x}}_i}{\partial \mathbf{x}_i} \times \frac{1}{\sqrt{\sigma_B^2 + \epsilon}} + \frac{\partial l}{\partial \sigma_B^2} \times \frac{2 \times (\mathbf{x}_i - \mu_B)}{m} + \frac{\partial l}{\partial \mu_B} \times \frac{1}{m}$$

$$⑤ \quad \frac{\partial l}{\partial \gamma} = \sum_{i=1}^m \frac{\partial l}{\partial \mathbf{y}_i} \times \hat{\mathbf{x}}_i$$

$$⑥ \quad \frac{\partial l}{\partial \beta} = \sum_{i=1}^m \frac{\partial l}{\partial \mathbf{y}_i}$$

Training Batch Normalization Networks

Input: Network N with trainable parameters Θ ; subset of activations $\{x^{(k)}\}_{k=1}^K$

Output: Batch-normalized network for inference N_{BN}^{inf}

- 1 $N_{BN}^{tr} = N$ // Training BN network
- 2 for $k = 1 \dots K$ do
- 3 Add transformation $y^{(k)} = BN_{\gamma^{(k)}, \beta^{(k)}}(x^{(k)})$ to N_{BN}^{tr}
- 4 Modify each layer in N_{BN}^{tr} with input $x^{(k)}$ to take $y^{(k)}$ instead
- 5 Train N_{BN}^{tr} to optimize the parameters $\Theta \cup \{\gamma^{(k)}, \beta^{(k)}\}_{k=1}^K$
- 6 $N_{BN}^{inf} = N_{BN}^{tr}$ // Inference BN network with frozen parameters
- 7 for $k = 1 \dots K$ do
- 8 Process multiple training mini-batches B , each of size m , and average over them
- 9 $E[x] = E_B[\mu_B]$ and $Var[x] = \frac{m}{m-1} E_B[\sigma_B^2]$
- 10 In N_{BN}^{inf} , replace the transform $y = BN_{\gamma, \beta}(x)$ with
- 11
$$y = \frac{\gamma}{\sqrt{Var[x] + \epsilon}} \times x + \left[\beta - \frac{\gamma E[x]}{\sqrt{Var[x] + \epsilon}} \right]$$

Training Batch Normalization Networks

Input: Network N with trainable parameters Θ ; subset of activations $\{x^{(k)}\}_{k=1}^K$

Output: Batch-normalized network for inference N_{BN}^{inf}

- 1 $N_{BN}^{tr} = N$ // Training BN network
 - for $k = 1 \dots K$ do
 - Add transformation $y^{(k)} = BN_{\gamma^{(k)}, \beta^{(k)}}(x^{(k)})$ to N_{BN}^{tr}
 - Modify each layer in N_{BN}^{tr} with input $x^{(k)}$ to take $y^{(k)}$ instead
 - Train N_{BN}^{tr} to optimize the parameters $\Theta \cup \{\gamma^{(k)}, \beta^{(k)}\}_{k=1}^K$
 - $N_{BN}^{inf} = N_{BN}^{tr}$ // Inference BN network with frozen parameters
 - for $k = 1 \dots K$ do
 - Process multiple training mini-batches B , each of size m , and average over them
 - $E[x] = E_B[\mu_B]$ and $Var[x] = \frac{m}{m-1} E_B[\sigma_B^2]$
 - In N_{BN}^{inf} , replace the transform $y = BN_{\gamma, \beta}(x)$ with
- $$y = \frac{\gamma}{\sqrt{Var[x] + \epsilon}} \times x + \left[\beta - \frac{\gamma E[x]}{\sqrt{Var[x] + \epsilon}} \right]$$

Training Batch Normalization Networks

Input: Network N with trainable parameters Θ ; subset of activations $\{x^{(k)}\}_{k=1}^K$

Output: Batch-normalized network for inference N_{BN}^{inf}

1 $N_{BN}^{tr} = N$ // Training BN network

2 for $k = 1 \dots K$ do

 Add transformation $y^{(k)} = BN_{\gamma^{(k)}, \beta^{(k)}}(x^{(k)})$ to N_{BN}^{tr}

 Modify each layer in N_{BN}^{tr} with input $x^{(k)}$ to take $y^{(k)}$ instead

 Train N_{BN}^{tr} to optimize the parameters $\Theta \cup \{\gamma^{(k)}, \beta^{(k)}\}_{k=1}^K$

$N_{BN}^{inf} = N_{BN}^{tr}$ // Inference BN network with frozen parameters

 for $k = 1 \dots K$ do

 Process multiple training mini-batches B , each of size m , and average over them

$E[x] = E_B[\mu_B]$ and $Var[x] = \frac{m}{m-1} E_B[\sigma_B^2]$

 In N_{BN}^{inf} , replace the transform $y = BN_{\gamma, \beta}(x)$ with

$$y = \frac{\gamma}{\sqrt{Var[x] + \epsilon}} \times x + \left[\beta - \frac{\gamma E[x]}{\sqrt{Var[x] + \epsilon}} \right]$$

Training Batch Normalization Networks

Input: Network N with trainable parameters Θ ; subset of activations $\{\mathbf{x}^{(k)}\}_{k=1}^K$

Output: Batch-normalized network for inference N_{BN}^{inf}

- 1 $N_{BN}^{tr} = N$ // Training BN network
- 2 for $k = 1 \dots K$ do
- 3 Add transformation $y^{(k)} = BN_{\gamma^{(k)}, \beta^{(k)}}(\mathbf{x}^{(k)})$ to N_{BN}^{tr}
 - 4 Modify each layer in N_{BN}^{tr} with input $\mathbf{x}^{(k)}$ to take $y^{(k)}$ instead
 - 5 Train N_{BN}^{tr} to optimize the parameters $\Theta \cup \{\gamma^{(k)}, \beta^{(k)}\}_{k=1}^K$
 - 6 $N_{BN}^{inf} = N_{BN}^{tr}$ // Inference BN network with frozen parameters
 - 7 for $k = 1 \dots K$ do
 - 8 Process multiple training mini-batches B , each of size m , and average over them
 - 9 $E[x] = E_B[\mu_B]$ and $Var[x] = \frac{m}{m-1} E_B[\sigma_B^2]$
 - 10 In N_{BN}^{inf} , replace the transform $y = BN_{\gamma, \beta}(x)$ with
 - 11
$$y = \frac{\gamma}{\sqrt{Var[x] + \epsilon}} \times x + \left[\beta - \frac{\gamma E[x]}{\sqrt{Var[x] + \epsilon}} \right]$$

Training Batch Normalization Networks

Input: Network N with trainable parameters Θ ; subset of activations $\{\mathbf{x}^{(k)}\}_{k=1}^K$

Output: Batch-normalized network for inference N_{BN}^{inf}

- 1 $N_{BN}^{tr} = N$ // Training BN network
- 2 for $k = 1 \dots K$ do
- 3 Add transformation $\mathbf{y}^{(k)} = BN_{\gamma^{(k)}, \beta^{(k)}}(\mathbf{x}^{(k)})$ to N_{BN}^{tr}
- 4 Modify each layer in N_{BN}^{tr} with input $\mathbf{x}^{(k)}$ to take $\mathbf{y}^{(k)}$ instead

5 Train N_{BN}^{tr} to optimize the parameters $\Theta \cup \{\gamma^{(k)}, \beta^{(k)}\}_{k=1}^K$

6 $N_{BN}^{inf} = N_{BN}^{tr}$ // Inference BN network with frozen parameters

7 for $k = 1 \dots K$ do

8 Process multiple training mini-batches B , each of size m , and average over them

9 $E[x] = E_B[\mu_B]$ and $Var[x] = \frac{m}{m-1} E_B[\sigma_B^2]$

10 In N_{BN}^{inf} , replace the transform $y = BN_{\gamma, \beta}(x)$ with

$$y = \frac{\gamma}{\sqrt{Var[x] + \epsilon}} \times x + \left[\beta - \frac{\gamma E[x]}{\sqrt{Var[x] + \epsilon}} \right]$$

Training Batch Normalization Networks

Input: Network N with trainable parameters Θ ; subset of activations $\{\mathbf{x}^{(k)}\}_{k=1}^K$

Output: Batch-normalized network for inference N_{BN}^{inf}

- 1 $N_{BN}^{tr} = N$ // Training BN network
- 2 for $k = 1 \dots K$ do
- 3 Add transformation $\mathbf{y}^{(k)} = BN_{\gamma^{(k)}, \beta^{(k)}}(\mathbf{x}^{(k)})$ to N_{BN}^{tr}
- 4 Modify each layer in N_{BN}^{tr} with input $\mathbf{x}^{(k)}$ to take $\mathbf{y}^{(k)}$ instead
- 5 Train N_{BN}^{tr} to optimize the parameters $\Theta \cup \{\gamma^{(k)}, \beta^{(k)}\}_{k=1}^K$

6 $N_{BN}^{inf} = N_{BN}^{tr}$ // Inference BN network with frozen parameters

7 for $k = 1 \dots K$ do

8 Process multiple training mini-batches B , each of size m , and average over them

9 $E[x] = E_B[\mu_B]$ and $Var[x] = \frac{m}{m-1} E_B[\sigma_B^2]$

10 In N_{BN}^{inf} , replace the transform $y = BN_{\gamma, \beta}(x)$ with

11
$$y = \frac{\gamma}{\sqrt{Var[x] + \epsilon}} \times x + \left[\beta - \frac{\gamma E[x]}{\sqrt{Var[x] + \epsilon}} \right]$$

Training Batch Normalization Networks

Input: Network N with trainable parameters Θ ; subset of activations $\{\mathbf{x}^{(k)}\}_{k=1}^K$

Output: Batch-normalized network for inference N_{BN}^{inf}

- 1 $N_{BN}^{tr} = N$ // Training BN network
- 2 for $k = 1 \dots K$ do
- 3 Add transformation $\mathbf{y}^{(k)} = BN_{\gamma^{(k)}, \beta^{(k)}}(\mathbf{x}^{(k)})$ to N_{BN}^{tr}
- 4 Modify each layer in N_{BN}^{tr} with input $\mathbf{x}^{(k)}$ to take $\mathbf{y}^{(k)}$ instead
- 5 Train N_{BN}^{tr} to optimize the parameters $\Theta \cup \{\gamma^{(k)}, \beta^{(k)}\}_{k=1}^K$
- 6 $N_{BN}^{inf} = N_{BN}^{tr}$ // Inference BN network with frozen parameters
- 7 for $k = 1 \dots K$ do
- 8 Process multiple training mini-batches B , each of size m , and average over them
- 9 $E[x] = E_B[\mu_B]$ and $Var[x] = \frac{m}{m-1} E_B[\sigma_B^2]$
- 10 In N_{BN}^{inf} , replace the transform $y = BN_{\gamma, \beta}(x)$ with
- 11
$$y = \frac{\gamma}{\sqrt{Var[x] + \epsilon}} \times x + \left[\beta - \frac{\gamma E[x]}{\sqrt{Var[x] + \epsilon}} \right]$$

Training Batch Normalization Networks

Input: Network N with trainable parameters Θ ; subset of activations $\{\mathbf{x}^{(k)}\}_{k=1}^K$

Output: Batch-normalized network for inference N_{BN}^{inf}

- 1 $N_{BN}^{tr} = N$ // Training BN network
- 2 for $k = 1 \dots K$ do
- 3 Add transformation $\mathbf{y}^{(k)} = BN_{\gamma^{(k)}, \beta^{(k)}}(\mathbf{x}^{(k)})$ to N_{BN}^{tr}
- 4 Modify each layer in N_{BN}^{tr} with input $\mathbf{x}^{(k)}$ to take $\mathbf{y}^{(k)}$ instead
- 5 Train N_{BN}^{tr} to optimize the parameters $\Theta \cup \{\gamma^{(k)}, \beta^{(k)}\}_{k=1}^K$
- 6 $N_{BN}^{inf} = N_{BN}^{tr}$ // Inference BN network with frozen parameters
- 7 for $k = 1 \dots K$ do

1 Process multiple training mini-batches B , each of size m , and average over them

2 $E[x] = E_B[\mu_B]$ and $Var[x] = \frac{m}{m-1} E_B[\sigma_B^2]$

3 In N_{BN}^{inf} , replace the transform $y = BN_{\gamma, \beta}(x)$ with

$$y = \frac{\gamma}{\sqrt{Var[x] + \epsilon}} \times x + \left[\beta - \frac{\gamma E[x]}{\sqrt{Var[x] + \epsilon}} \right]$$

Training Batch Normalization Networks

Input: Network N with trainable parameters Θ ; subset of activations $\{\mathbf{x}^{(k)}\}_{k=1}^K$

Output: Batch-normalized network for inference N_{BN}^{inf}

- 1 $N_{BN}^{tr} = N$ // Training BN network
- 2 for $k = 1 \dots K$ do
- 3 Add transformation $\mathbf{y}^{(k)} = BN_{\gamma^{(k)}, \beta^{(k)}}(\mathbf{x}^{(k)})$ to N_{BN}^{tr}
- 4 Modify each layer in N_{BN}^{tr} with input $\mathbf{x}^{(k)}$ to take $\mathbf{y}^{(k)}$ instead
- 5 Train N_{BN}^{tr} to optimize the parameters $\Theta \cup \{\gamma^{(k)}, \beta^{(k)}\}_{k=1}^K$
- 6 $N_{BN}^{inf} = N_{BN}^{tr}$ // Inference BN network with frozen parameters
- 7 for $k = 1 \dots K$ do
- 8 Process multiple training mini-batches \mathcal{B} , each of size m , and average over them



$$E[x] = E_{\mathcal{B}}[\mu_{\mathcal{B}}] \text{ and } Var[x] = \frac{1}{m-1} E_{\mathcal{B}}[\sigma_{\mathcal{B}}^2]$$



In N_{BN}^{inf} , replace the transform $y = BN_{\gamma, \beta}(x)$ with



$$y = \frac{\gamma}{\sqrt{Var[x] + \epsilon}} \times x + \left[\beta - \frac{\gamma E[x]}{\sqrt{Var[x] + \epsilon}} \right]$$

Training Batch Normalization Networks

Input: Network N with trainable parameters Θ ; subset of activations $\{\mathbf{x}^{(k)}\}_{k=1}^K$

Output: Batch-normalized network for inference N_{BN}^{inf}

- 1 $N_{BN}^{tr} = N$ // Training BN network
- 2 for $k = 1 \dots K$ do
- 3 Add transformation $\mathbf{y}^{(k)} = BN_{\gamma^{(k)}, \beta^{(k)}}(\mathbf{x}^{(k)})$ to N_{BN}^{tr}
- 4 Modify each layer in N_{BN}^{tr} with input $\mathbf{x}^{(k)}$ to take $\mathbf{y}^{(k)}$ instead
- 5 Train N_{BN}^{tr} to optimize the parameters $\Theta \cup \{\gamma^{(k)}, \beta^{(k)}\}_{k=1}^K$
- 6 $N_{BN}^{inf} = N_{BN}^{tr}$ // Inference BN network with frozen parameters
- 7 for $k = 1 \dots K$ do
- 8 Process multiple training mini-batches \mathcal{B} , each of size m , and average over them
- 9 $E[\mathbf{x}] = E_{\mathcal{B}}[\mu_{\mathcal{B}}]$ and $Var[\mathbf{x}] = \frac{m}{m-1} E[\sigma_{\mathcal{B}}^2]$

10 In N_{BN}^{inf} , replace the transform $\mathbf{y} = BN_{\gamma, \beta}(\mathbf{x})$ with

$$\mathbf{y} = \frac{\gamma}{\sqrt{Var[\mathbf{x}] + \epsilon}} \times \mathbf{x} + \left[\beta - \frac{\gamma E[\mathbf{x}]}{\sqrt{Var[\mathbf{x}] + \epsilon}} \right]$$

Training Batch Normalization Networks

Input: Network N with trainable parameters Θ ; subset of activations $\{\mathbf{x}^{(k)}\}_{k=1}^K$

Output: Batch-normalized network for inference N_{BN}^{inf}

- 1 $N_{BN}^{tr} = N$ // Training BN network
- 2 for $k = 1 \dots K$ do
- 3 Add transformation $\mathbf{y}^{(k)} = BN_{\gamma^{(k)}, \beta^{(k)}}(\mathbf{x}^{(k)})$ to N_{BN}^{tr}
- 4 Modify each layer in N_{BN}^{tr} with input $\mathbf{x}^{(k)}$ to take $\mathbf{y}^{(k)}$ instead
- 5 Train N_{BN}^{tr} to optimize the parameters $\Theta \cup \{\gamma^{(k)}, \beta^{(k)}\}_{k=1}^K$
- 6 $N_{BN}^{inf} = N_{BN}^{tr}$ // Inference BN network with frozen parameters
- 7 for $k = 1 \dots K$ do
- 8 Process multiple training mini-batches \mathcal{B} , each of size m , and average over them
- 9 $E[\mathbf{x}] = E_{\mathcal{B}}[\mu_{\mathcal{B}}]$ and $Var[\mathbf{x}] = \frac{m}{m-1} E[\sigma_{\mathcal{B}}^2]$
- 10 In N_{BN}^{inf} , replace the transform $\mathbf{y} = BN_{\gamma, \beta}(\mathbf{x})$ with
- 11
$$\mathbf{y} = \frac{\gamma}{\sqrt{Var[\mathbf{x}] + \epsilon}} \times \mathbf{x} + \left[\beta - \frac{\gamma E[\mathbf{x}]}{\sqrt{Var[\mathbf{x}] + \epsilon}} \right]$$

However

Santurkar et al. [16]

- They found that it is not the covariance shift that is affected by it!!!

Santurkar et al. recognize that

- Batch normalization has been arguably one of the most successful architectural innovations in deep learning.

They used a standard very deep convolutional network

- on CIFAR-10 with and without BatchNorm

However

Santurkar et al. [16]

- They found that it is not the covariance shift that is affected by it!!!

Santurkar et al. recognize that

- Batch normalization has been arguably one of the most successful architectural innovations in deep learning.

They used a standard vanilla deep convolutional neural network

- on CIFAR-10 with and without BatchNorm

However

Santurkar et al. [16]

- They found that is not the covariance shift the one affected by it!!!

Santurkar et al. recognize that

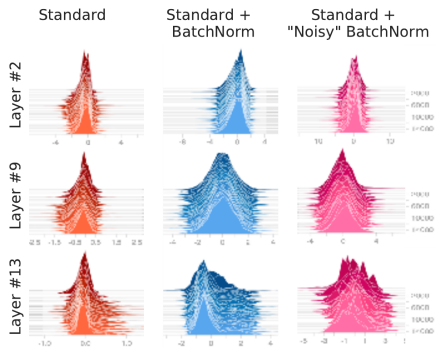
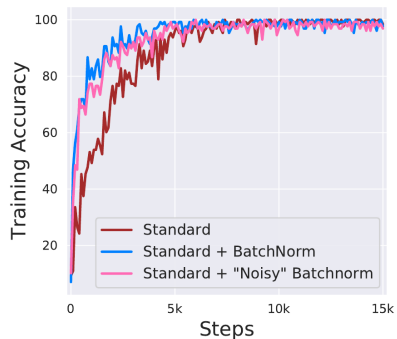
- Batch normalization has been arguably one of the most successful architectural innovations in deep learning.

They used a standard Very deep convolutional network

- on CIFAR-10 with and without BatchNorm

They found something quite interesting

The following facts



Actually Batch Normalization

It does not do anything to the Internal Covariate Shift

- Actually smooth the optimization manifold
 - ▶ It is not the only way to achieve it!!!

They suggest that:

- “This suggests that the positive impact of BatchNorm on training might be somewhat serendipitous.”

Actually Batch Normalization

It does not do anything to the Internal Covariate Shift

- Actually smooth the optimization manifold
 - ▶ It is not the only way to achieve it!!!

They suggest that

- “This suggests that the positive impact of BatchNorm on training might be somewhat serendipitous.”

They actually have a connected result

To the analysis of gradient clipping!!!

- They are the same group

Theorem (The effect of BatchNorm on the Lipschitzness of the loss)

- For a BatchNorm network with loss $\hat{\mathcal{L}}$ and an identical non-BN network with (identical) loss \mathcal{L} ,

$$\|\nabla_{y_j} \hat{\mathcal{L}}\|^2 \leq \frac{\gamma^2}{\sigma_j^2} \left[\|\nabla_{y_j} \mathcal{L}\|^2 - \frac{1}{m} \langle \mathbf{1}, \nabla_{y_j} \mathcal{L} \rangle^2 - \frac{1}{\sqrt{m}} \langle \nabla_{y_j} \mathcal{L}, \hat{y}_j \rangle^2 \right]$$

They actually have a connected result

To the analysis of gradient clipping!!!

- They are the same group

Theorem (The effect of BatchNorm on the Lipschitzness of the loss)

- For a BatchNorm network with loss $\hat{\mathcal{L}}$ and an identical non-BN network with (identical) loss \mathcal{L} ,

$$\|\nabla_{\mathbf{y}_j} \hat{\mathcal{L}}\|^2 \leq \frac{\gamma^2}{\sigma_j^2} \left[\|\nabla_{\mathbf{y}_j} \mathcal{L}\|^2 - \frac{1}{m} \langle \mathbf{1}, \nabla_{\mathbf{y}_j} \mathcal{L} \rangle^2 - \frac{1}{\sqrt{m}} \langle \nabla_{\mathbf{y}_j} \mathcal{L}, \hat{\mathbf{y}}_j \rangle^2 \right]$$

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
 - Deconvolution Layer
 - Alternating Minimization
- 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
- Normalization Layer AKA Batch Normalization
- **Finally, The Fully Connected Layer**

4 An Example of CNN

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

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.

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

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

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 .

Something to Note

- 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.

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

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

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

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

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

Basically

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

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

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

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

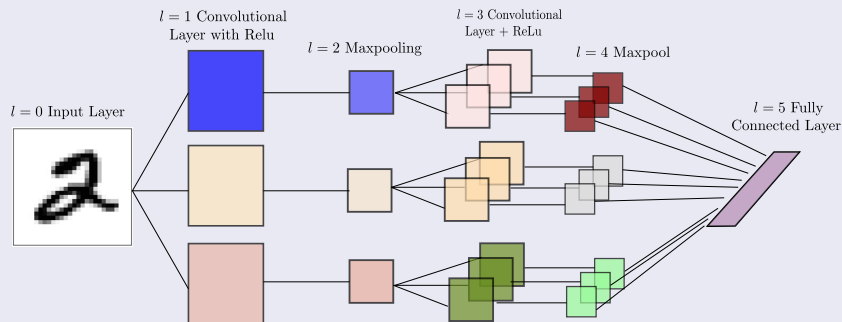
- Convolutional Layer
 - Convolutional Architectures
 - A Little Bit of Notation
 - Deconvolution Layer
 - Alternating Minimization
- Non-Linearity Layer
 - Fixing the Problem, ReLu function
 - Back to the Non-Linearity Layer
- Rectification Layer
- Local Contrast Normalization Layer
- Sub-sampling and Pooling
 - Strides
- Normalization Layer AKA Batch Normalization
- 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"



Therefore, we have

Layer $l = 1$

- This Layer is using a ReLu f with 3 channels

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

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

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

Layer $l = 2$

We have a maxpooling of size 2×2

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

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

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

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
 - Deconvolution Layer
 - Alternating Minimization
- 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
- Normalization Layer AKA Batch Normalization
- 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):

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

Therefore, we can obtain

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

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

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

With

$$y_1^{(6)} = \text{ReLU}(z_1^{(5)})$$

Therefore

We get in the first part of the equation

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

With

$$y_1^{(6)} = \text{ReLU}(z_1^{(5)})$$

Therefore

We have

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

Therefore, if we use the approximation

$$\frac{\partial f(z_1^{(5)})}{\partial z_1^{(5)}} = \frac{e^{kz_1^{(5)}}}{(1 + e^{kz_1^{(5)}})^2}$$

Therefore

We have

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

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

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

Finally

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

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

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

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 $j = 1$

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

Therefore

We have then

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

We have the following chain of derivations given

$$[X^{(6)}] = f [Y^{(4)}]$$

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

Therefore

We have then

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

We have the following chain of derivations given

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

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

Therefore

We have

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

Then assuming that

$$(Y_f^{(3)})_{x,y} = (B_f^{(3)})_{x,y} + \sum_{k=-h_1^{(0)}}^{h_1^{(0)}} \sum_{l=-h_2^{(0)}}^{h_2^{(0)}} (K_{f_1}^{(3)})_{k,l} (Y_1^{(2)})_{x-k,y-l}$$

Therefore

We have

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

Then assuming that

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

Therefore

We have

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

Then

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

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 Other Equations

I will leave you to devise them

- They are a repetitive procedure.

The interesting one is the average pooling

- The others are the stride and the deconvolution

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. Zahoora, 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] M. D. Zeiler, D. Krishnan, G. W. Taylor, and R. Fergus, "Deconvolutional networks," in *2010 IEEE Computer Society Conference on computer vision and pattern recognition*, pp. 2528–2535, IEEE, 2010.
- [9] D. Krishnan and R. Fergus, "Fast image deconvolution using hyper-laplacian priors," *Advances in neural information processing systems*, vol. 22, pp. 1033–1041, 2009.

- [10] Y. Wang, J. Yang, W. Yin, and Y. Zhang, "A new alternating minimization algorithm for total variation image reconstruction," *SIAM Journal on Imaging Sciences*, vol. 1, no. 3, pp. 248–272, 2008.
- [11] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. The MIT Press, 2016.
- [12] J. T. Springenberg, A. Dosovitskiy, T. Brox, and M. Riedmiller, "Striving for simplicity: The all convolutional net," 2015.
- [13] T. G. Kolda and B. W. Bader, "Tensor decompositions and applications," *SIAM review*, vol. 51, no. 3, pp. 455–500, 2009.
- [14] C. Kong and S. Lucey, "Take it in your stride: Do we need striding in cnns?," *arXiv preprint arXiv:1712.02502*, 2017.
- [15] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," *arXiv preprint arXiv:1502.03167*, 2015.

- [16] S. Santurkar, D. Tsipras, A. Ilyas, and A. Madry, “How does batch normalization help optimization?,” in *Advances in Neural Information Processing Systems*, pp. 2483–2493, 2018.