

Universität Hamburg
Department Informatik
Knowledge Technology, WTM

Deep Learning: Neural Networks for Object Detection and Tracking Tasks

Seminar Paper

Brain Modelling

Daniel Speck

Matr.Nr. 632 13 17

2speck@informatik.uni-hamburg.de

29.05.2015

Abstract

Deep neural networks are one of the most successful learning strategies at the moment as the computing power for creating such structures rised in the past years via GPU computing. Object detection and tracking tasks can be fulfilled with these architectures.

Contents

1	Introduction	2
2	Background information: Artificial neural networks	2
3	Deep neural networks	2
4	Convolutional neural networks and image processing	2
5	Research and field of application	5
5.1	Image classification	5
5.2	Video classification	5
5.3	Object tracking	5
6	Conclusion	5
	Bibliography	6

1 Introduction

Deep learning is subcategory of machine learning and the focus of this paper will be deep neural networks in the context of deep learning.

An overview of image classification will be made [1] [3].

The visual cortex and deep learning strategies will be introduced [4].

Approaches for object detection [5] and tracking [6] via deep neural networks will be discussed.

2 Background information: Artificial neural networks

Artificial neural networks are intended to approximate certain functions for machine learning purposes. Fix/static algorithms can be calculated fast by modern computers but fail most times at disciplines requiring "intelligent" behavior. A classic example is recognizing handwritten digits because the shape, color and contrast highly vary in dependence of the used pencil and, of course, of the writing style.

3 Deep neural networks

Classic, small neural networks can not succeed on complex tasks such as object detection in high-resolution, real world images. The variety of features, blurry backgrounds, one or several objects in an image and other factors render detection, tracking and classification tasks challenging. The ability to detect, track and distinguish real world objects demands complex structures with a much larger capability of solving complex problems.

As computing power rises and GPU-computing became popular neural networks are no longer limited to a few layers containing just some neurons. Modern solutions (2014) can carry tens of layers with thousands of neurons and millions of connections [2]. Complex structures like this enable computing complex problems which arise with tasks like image classification in real world images.

4 Convolutional neural networks and image processing

For deep learning purposes (classic / fully-connected) multilayer perceptrons consume a sizable amount of resources for proper training when they are designed to solve complex tasks because the amount of neurons and especially weights increases rapidly with the network's size.

For example, a MLP with three layers, an input layer with 100 neurons, a hidden layer with 25 neurons and an output layer with 10 neurons for classifying images

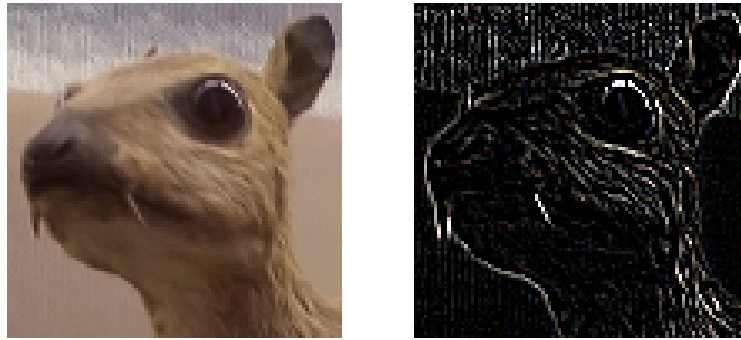


Figure 1: Left side: original image, right side: edge-detection kernel processed image. Original image by Michael Plotke, 28th of January, 2013. Open creative commons license. <http://upload.wikimedia.org/wikipedia/commons/5/50/Vd-Orig.png> and <http://upload.wikimedia.org/wikipedia/commons/6/6d/Vd-Edge3.png>

with a size of 10x10 pixels into 10 different classes would have $100 * 25 + 25 * 10 = 2,750$ weights/connections. Training this net would already result in a big time and space complexity. Moreover, as features in images capturing real world scenes are distributed in certain patterns (they cover spatially local correlation, such as shapes), there is no need to have every pixels information being processed by one neuron. Actually in most cases results would be even better, if the pixels information is pre-processed, for instance by edge detection filters but a fully connected layer of neurons is not an optimal solution for this task.

Convolutional neural networks (CNNs) are inspired by biology, instead of connecting every pixels information directly with a neuron to process its information it filters the information in the first layers [3]. This procedure is similar to the on processes happening when an biological eye receives stimuli.

The receptive field ¹ has a vast amount of photoreceptor cells ² gathering information and converging the received information on to distinctly less retinal ganglion cells ³. This process maps several features and reduces the input dimensionality as well as distinguishes the information to separate "channels" which are then transferred to the corresponding neurons to process features such as color, motion, shapes and so on separately [4]. The idea of CNNs is based on this biological processes, the information of an input image is convolved by several filters which try to extract interesting features in the first layer and in following layers this information is pooled and subsampled [3].

Convolution itself is the applying of a function repeatedly of the output of another function and in the context of CNNs it is applying different "filters" over an image to extract the already mentioned features. A convolution layer extracts the pixel information out of an image with kernels ⁴.

¹http://en.wikipedia.org/wiki/Receptive_field

²http://en.wikipedia.org/wiki/Photoreceptor_cell

³http://en.wikipedia.org/wiki/Retinal_ganglion_cell

⁴[http://en.wikipedia.org/wiki/Kernel_\(image_processing\)](http://en.wikipedia.org/wiki/Kernel_(image_processing))

Example for a convolution layer: In figure 1 you can see an image of an animal on the left side (original image) and a kernel processed one on the right side. The used kernel matrix for filtering the left image is shown in equation (1). So basically each pixels information in the right, processed image is the result of applying the kernel matrix (1) on the same pixel (and the neighboring pixels) in the left image. With a wide variety of different kernels several different features can be extracted from an image.

$$K_M = \begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix} \quad (1)$$

Typical CNNs use tens to hundreds of different kernel filters gathering as much features from an image as possible. The resulting pixel processed by a kernel filter is calculated via the formula:

$$I_{out}(x, y) = \sum_{a=1}^3 \sum_{b=1}^3 I_{in}(x + a - c_x, y + b - c_y) * K_M(a, b) \quad (2)$$

Where c_x is the coordinate of the x-center and c_y the coordinate of the y-center of the input image. For some example input image I_{in} with a dimension of 3x3 pixels the center pixel (coordinates $x = 2, y = 2$) of the kernel processed output image I_{out} using the kernel filter K_M of equation (1) would be calculated like this:

$$I_{in} = \begin{bmatrix} 46 & 42 & 50 \\ 44 & 65 & 56 \\ 41 & 52 & 58 \end{bmatrix}$$

$$\begin{aligned} I_{out}(2, 2) &= 46 * (-1) + 42 * (-1) + 50 * (-1) + 44 * (-1) + 55 * 8 + 56 * (-1) \\ &\quad + 41 * (-1) + 52 * (-1) + 58 * (-1) = 131 \end{aligned}$$

After processing this filter to the whole image edges would be highlighted and the rest nearly black, like in figure 1 so that shape features are extracted out of the original image.

Another idea of CNNs is subsampling layers following convolutional layers. Subsampling reduces the overall size of information and therefore not only saves resources for the later classification tasks but also strengthens the detected features in an image. Often max-pooling is used as a subsampling strategy in CNNs [5, 3] which determines the most distinctive pixel in a given area. A max-pooling algorithm splits the input in grids and selects the maximum value out of each grid, thus non-maximal values are deleted so that only those information continues to exist which best represents the current feature. Before applying max-pooling the input represents the presence of a feature in one or only some pixels. After this kind of dimensionality reduction the assertion is enlarged to whole area, corresponding to

several pixels of the original image.

Additionally this technique provides robustness to the position of a feature in an image, as the position becomes less important.

The last layers of a CNN regularly consist of fully-connected layers. In comparison to feature extraction, strengthening and dimensionality reduction those layers are supposed to take all gathered features / information and make predictions. If the supplied input contains big shapes of ears, legs etc., maybe filtered trunks, grey-ish color maps and so on the fully-connected layers have to collect and mix all those gathered features / information and to recognize it for finally classifying it as an image of elephants.

5 Research and field of application

5.1 Image classification

More details / information about state of the art object/image detection/classification.

5.2 Video classification

More details / information about state of the art object/video detection/classification.

5.3 Object tracking

6 Conclusion

Conclusion of the paper.

References

- [1] D. Ciresan, U. Meier, and J. Schmidhuber. Multi-column deep neural networks for image classification. *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*, pages 3642 – 3649, June 2012.
- [2] Andrej Karpathy, George Toderici, Sanketh Shetty, Thomas Leung, Rahul Sukthankar, and Li Fei-Fei. Large-scale video classification with convolutional neural networks. In *CVPR*, 2014.
- [3] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep convolutional neural networks. In F. Pereira, C.J.C. Burges, L. Bottou, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 1097–1105. Curran Associates, Inc., 2012.
- [4] Norbert Kruger, Peter Janssen, Sinan Kalkan, Markus Lappe, Ales Leonardis, Justus Piater, Antonio Jose Rodriguez-Sanchez, and Laurenz Wiskott. Deep hierarchies in the primate visual cortex: What can we learn for computer vision? *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 35(8):1847–1871, 2013.
- [5] Christian Szegedy, Alexander Toshev, and Dumitru Erhan. Deep neural networks for object detection. In C.J.C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems 26*, pages 2553–2561. Curran Associates, Inc., 2013.
- [6] Naiyan Wang and Dit-Yan Yeung. Learning a deep compact image representation for visual tracking. In C.J.C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems 26*, pages 809–817. Curran Associates, Inc., 2013.