

Object tracking with deep neural networks

Santeri Salmijärvi

School of Electrical Engineering

Bachelor's thesis

Espoo Work in progress! Compiled: 17:30:50 2017/07/07

Thesis supervisor:

D.Sc. (Tech) Pekka Forsman

Thesis advisor:

M.Sc. (Tech) Mikko Vihlman

Author: Santeri Salmijärvi

Title: Object tracking with deep neural networks

Date: Work in progress! Compiled: 17:30:50 2017/07/07 Language: English
Number of pages: 5+8

Degree programme: Bachelor's Program in Electrical Engineering

Supervisor: D.Sc. (Tech) Pekka Forsman

Advisor: M.Sc. (Tech) Mikko Vihlman

abstract in english

Keywords: keywords in english

Tekijä: Santeri Salmijärvi

Työn nimi: Kohteenseuranta syvillä neuroverkoilla

Päivämäärä: Work in progress! Compiled: 17:30:50 2017/07/07 Kieli: Englanti
Sivumäärä: 5+8

Koulutusohjelma: Sähkötekniikan kandidaattiohjelma

Vastuuopettaja: TkT Pekka Forsman

Työn ohjaaja: DI Mikko Vihlman

lyhyt tiivistelmä suomeksi

Avainsanat: avainsanat suomeksi

Contents

Abstract	ii
Abstract (in Finnish)	iii
Contents	iv
1 Introduction	1
2 Deep Learning	2
2.1 Deep neural networks	2
2.2 Convolutional networks	2
2.3 Stacked denoising autoencoders	3
3 Object tracking	4
3.1 Target representation	4
3.2 Model update	4
3.3 Challenges	4
4 Deep neural networks in tracking	5
5 Data sets and evaluation	6
6 Conclusions	7
References	8

Abbreviations

CNN Convolutional Neural Network

DNN Deep Neural Network

MLP MultiLayer Perceptron

NN Neural Network

ReLU Rectified Linear Unit

SDAE Stacked Denoising Autoencoder

1 Introduction

Object tracking is a large and actively researched sub-area of computer vision. The main task for a tracker is to find and follow the desired subject in a sequence of images. Object tracking is closely related to other image analysis tasks so the implementations also share elements. In the recent years, use of deep neural networks has been researched for object tracking.

Many of the deep networks tailored to tracking tasks are variations of convolutional networks. Another way used to extract features from a frame is a stacked denoising autoencoder [1]. The training of deep neural networks requires a large amount of training data and their development has been made easier by an increase in the size of appli-

capable datasets.

Comment by author:
add cites to examples

The goal of this thesis is to study the concepts behind object tracking and deep neural networks. It will present the architectures and principles currently used in deep neural networks tailored to object tracking tasks. The practices and datasets used in training and evaluating such networks are also introduced.

2 Deep Learning

This chapter introduces the basic concepts behind deep neural networks and the first two sub-sections are based on the book Deep Learning by Goodfellow et. al. [2]

2.1 Deep neural networks

A Deep Neural Network (DNN) is commonly defined as a Neural Network (NN), that has a **visible** input and output layer with several **hidden layers** between them. The distinction between visible and hidden layers is important because training of the network only evaluates the output layer's performance. During training, a **learning algorithm** optimizes the individual hidden layers to best approximate the desired output of the whole network.

The input layer takes in the data to be processed, which typically means a vector of color values in the case of object tracking. These are then processed by the hidden layers and finally the output layer produces the target's position in the frame. These models usually come in the form of a **feedforward neural network** or **MultiLayer Perceptron (MLP)**. The name comes from the fact that information flows from the input through computations to the output with no **feedback** connections.

In NNs, each layer consist of several **units** with a weight and activation function. A bias-term can also be defined for each unit. The weights of a layer are commonly represented by a matrix by which the input-vector is multiplied. Units in a layer also have a common activation function, which is fed by the sum of its weighted inputs in addition to the possible bias, and the result is output to the next layer alongside the layers other units' outputs. A commonly used unit type is the **Rectified Linear Unit (ReLU)**, which is defined by the activation function $g(z) = \max\{0, z\}$. It provides a nonlinear transformation while being comparable to linear models in terms of generalizing well and being easy to optimize.

Comment by author:

picture from eg. deeplearningbook page 174?

Before training, the weights of a MLP are initialized to small random values and biases to zero or small positive values. Then an algorithm called **stochastic gradient descent** is commonly applied alongside a training dataset. The basic procedure is to calculate the error of the network's output values compared to the desired ones using a **loss function**. The function's gradient can then be calculated for example by **back-propagation**, which feeds the errors back through the network to assign a contribution value to each unit. These values are then used to calculate the gradient of the loss function relative to the weights. Each weight is adjusted slightly to the opposite sign to minimize the loss function.

2.2 Convolutional networks

A Convolutional Neural Network (CNN) is simply a NN that uses convolution instead of general matrix multiplication in at least one of its layers. The main benefits of convolution in NNs are that it's dramatically more efficient in terms of memory

requirements, it reduces the amount of computation needed and it makes it possible to work with variable input sizes.

A typical convolutional layer consists of three stages: a convolution stage, detector stage and pooling stage. These can be implemented by individual layers. First, a **kernel** is applied to the input data in positions separated by a stepsize. This means that a linear activation function is fed by the matrix product of the input location and the kernel's weight matrix. In the detector stage, the results are then run through a non-linear activation, for example a ReLU. Finally, a **pooling function** is used to combine the results of multiple nearby outputs as the final output.

Comment by author:
pictures from dl p.337?

Convolutional layers enable indirect connections to all or most of the input data deeper in the network even when individual layers' connections are very sparse.

2.3 Stacked denoising autoencoders

A Stacked Denoising Autoencoder (SDAE) is a modification of the classic autoencoders. Denoising autoencoders are trained to encode a corrupted version of the input to a hidden representation and decode that to useful features of the clean input. A stacked denoising autoencoder is simply a sequence of denoising autoencoders trained this way. Corrupted input is only used to train the individual layers to find useful features so a trained SDAE works on clean input. [3]

Comment by author:
motivation for using sdaes?

3 Object tracking

Object tracking in video sequences has been researched for decades using different approaches for defining the target and adapting to changes in its shape or orientation. The situations most likely to cause tracking failure have also been identified.

3.1 Target representation

Tracking methods can be roughly divided to generative and discriminative, but combinations of them have also been proposed.

Generative methods search the frame for the best matches to a template of an appearance model of the subject. Template methods based on pixel intensity and color histograms perform well with no drastic changes in object appearance and non-cluttered backgrounds. Appearance models learned from training can be less affected by appearance variations and adaptive schemes provide added flexibility, while sparse models handle occlusion and image noise better. [4]

Discriminative methods consider tracking as a binary classification problem. They take the background also into account to separate the target from it. Used approaches include refining the initial guess with a support vector machine [5] or utilizing a relevance vector machine [6].

3.2 Model update

(Describe roughly the idea of updating the model online, motivations)

3.3 Challenges

(Present challenging situations for trackers, motivation for development of better methods.)

4 Deep neural networks in tracking

(Overview of task from DNN-point of view, strengths and weaknesses compared to more traditional solutions.)

An early implementation of a CNN-based tracker [7] pre-dates the work of Krizhevsky et. al. [8]. It takes modifies the architecture used for detection to make the network less affected by shifts in the objects position in the frame. Shift-invariancy is a non-desirable quality in tracking while using previous positions as a as it might result to mixups with objects similar to the target. [7]

5 Data sets and evaluation

(Overview of the data sets used for training and analysis. Methods used for comparing performance.)

The datasets used for training are equally important as the actual network design. Research on networks working with image data has been made easier by larger sets of both hand-labeled sets and ones obtained by simple keyword searches from online image services. These kinds of sets can be used to pre-train useful target features to tracking networks.

There has also been an increase in resources devoted to tracking data with the TR-100 -set [9] introduced in being a good example. It contains a hundred tracking sequences with reference positions for the target on each frame. Because some of the targets are similar or less challenging, a subset of 50 sequences considered challenging is also provided as TR-50. [10]

6 Conclusions

Summarize the current state of object tracking with DNNs with possibly some insight to future developments.

References

- [1] Wang, Naiyan and Dit Yan Yeung: *Learning a deep compact image representation for visual tracking*. In *Advances in Neural Information Processing Systems 26*, pages 809–817. Curran Associates, Inc., 2013, ISBN 9781632660244.
- [2] Goodfellow, I., Y. Bengio, and A. Courville: *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- [3] Vincent, P., H. Larochelle, I. Lajoie, Y. Bengio, and P. A. Manzagol: *Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion*. *Journal of Machine Learning Research*, 11:3371–3408, 2010. <http://www.jmlr.org/papers/v11/vincent10a.html>.
- [4] Wang, Q., F. Chen, W. Xu, and M. H Yang: *Object tracking via partial least squares analysis*. *IEEE Transactions on Image Processing*, 21(10):4454–4465, 2012.
- [5] Avidan, S.: *Support vector tracking*. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 26(8):1064–1072, 2004.
- [6] Williams, O., A. Blake, and R. Cipolla: *Sparse bayesian learning for efficient visual tracking*. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 27(8):1292–1304, 2005.
- [7] Fan, J., W. Xu, Y. Wu, and Y. Gong: *Human tracking using convolutional neural networks*. *IEEE Transactions on Neural Networks*, 21(10):1610–1623, 2010.
- [8] Krizhevsky, A., I. Sutskever, and G. E. Hinton: *Imagenet classification with deep convolutional neural networks*. In *Advances in Neural Information Processing Systems*, volume 2, pages 1097–1105, 2012. <http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf>.
- [9] *Visual tracker benchmark*. http://cvlab.hanyang.ac.kr/tracker_benchmark/index.html, Accessed: 2017-06-25.
- [10] Wu, Y., J. Lim, and M. H Yang: *Object tracking benchmark*. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(9):1834–1848, 2015.