

University of Warsaw
Faculty of Mathematics, Informatics and
Mechanics

**Krystyna Gajczyk, Jakub Pierewoj,
Przemysław Przybyszewski, Adam Starak**

Student number: 332118, 360641, 332493, 361021

Inference in neural networks using low-precision arithmetic

Bachelor thesis in COMPUTER SCIENCE

Supervised by:
Dr. Konrad Durnoga

June 17, 2017

Oświadczenie kierującego pracą

Potwierdzam, że niniejsza praca została przygotowana pod moim kierunkiem i kwalifikuje się do przedstawienia jej w postępowaniu o nadanie tytułu zawodowego.

Data

Podpis kierującego pracą

Oświadczenie autora (autorów) pracy

Świadom odpowiedzialności prawnej oświadczam, że niniejsza praca dyplomowa została napisana przeze mnie samodzielnie i nie zawiera treści uzyskanych w sposób niezgodny z obowiązującymi przepisami.

Oświadczam również, że przedstawiona praca nie była wcześniej przedmiotem procedur związanych z uzyskaniem tytułu zawodowego w wyższej uczelni.

Oświadczam ponadto, że niniejsza wersja pracy jest identyczna z załączoną wersją elektroniczną.

Data

Podpis autora (autorów) pracy

Abstract

TODO

Keywords

binarized neural network, XORNET

Subject Area (Socrates/Erasmus code)

11.3 Informatics, Computer Science

Categories and Subject Descriptors

Computing methodologies - Machine learning - Machine learning approaches - Neural Networks

Thesis title in Polish

Inferencja w sieciach neuronowych przy użyciu arytmetyki niskiej precyzji

Contents

1. Introduction	5
1.1. What is deep learning and why is it interesting?	5
1.2. Why binarized networks are interesting?	5
1.3. Goal of the project	6
2. Architecture overview	9
2.1. Used framework	9
2.2. Our implementation	9
2.3. Used networks and datasets	9
2.3.1. LeNet on MNIST	9
2.3.2. AlexNet on ImageNet	11
2.3.3. Finetuned AlexNet on oxford-102	11
2.3.4. Residual network on CIFAR-10	12
2.4. Implemented binarization algorithm	12
2.4.1. Binarized convolution filters from XORNET [2]	12
2.4.2. Binarized convolution filters inspired by DoReFa [4]	13
2.4.3. Binarized filters and activations	13
3. Experiments results	15
3.1. Lenet on MNIST	15
3.2. AlexNet on Flowers	15
3.3. Residual network on CIFAR-10	15
3.3.1. Results and discussion	15
4. Conclusions	17
4.1. Discussion of experiments results	17
4.2. Future work	17
5. Team members contribution	19
Bibliography	21

Chapter 1

Introduction

1.1. What is deep learning and why is it interesting?

Deep learning is a branch of machine learning which tries to model high level abstractions in data (like understanding words or finding objects on image) using a graph of simple elements - neurons, connected with specific activations and weights to others. It has recently gained a lot of interest from the industry, especially after recent successes in image or speech recognition.

The most important part of deep learning is training net based on the samples that the result is known. After each training step one can change weights and activation in model to make output of the network as close as possible to the expected value.

TODO:MIEJSCE NA SŁOWNICZEK i opis co to konwolucja

Currently deep neural networks get better results than state-of-the-art algorithms in many areas, such as computer vision and speech recognition. This breakthrough was possible because computing power of modern computers is high enough to handle intensive workload required to train large artificial neural networks.

1.2. Why binarized networks are interesting?

Deep neural networks (DNN) usually use high precision (floating point) numbers to represent weights and activations. Computations using that arithmetic are usually handled by graphic cards. On devices with low computing power, such as mobile phones, use of deep neural networks is very limited.

TODO:opis ile można zaoszczędzić przy różnych typach binaryzacji + inne potencjalne zyski

Researchers tried to address this problem in the number of papers in the recent years. A couple of main approaches have been investigated in order to reduce inefficient computation and memory usage in deep neural networks while maintaining the classification accuracy. Those approaches can be summarized as:

- Shallow networks - it has been proved, that every neural network

can be replaced with a corresponding neural network with a single (possibly large) hidden layer. One of the main problems with this approach is, that in order to achieve similar accuracy as the original neural network, the shallow networks need approximately the same number of parameters (numbers of neurons and connection between them). Second problem is that empirical test have shown, that while it shows good results for relatively small datasets, it underperforms for larger datasets (such as ImageNet).

- Compression of pre-trained DNN - it is possible to prune a DNN at inference time by discarding the weights, that don't bring any information to the classification process. It has been also further proposed to reduce the number of parameters and/or activations in order to increase acceleration, compression. This reduces the main blockers for using DNN on small, embedded devices, such as memory usage and energy. Memory usage can be achieved through quantizing and compressing weights with Huffman coding or hashing weights, such that the weights assigned to one hash bucket share the same parameter.
- Compact layers - this approach also reduces the memory usage and computation time. This approach replaces parts of the DNN structure with corresponding elements, which are smaller in size and bring nearly as much information. A few techniques have been examined, such as replacing a fully connected layer with global average pooling and replacing a convolution with a corresponding one requiring smaller amount of parameters.
- Quantizing parameters - this technique aims to replace the floating-point parameters of the neural network with the quantized values (through vector quantization methods), which require smaller number of bits of memory and need simpler arithmetic operations for computation. Number of DNN with quantization have been designed: 8-bit integer instead of 32-bit floating point activations, ternary weights and 3-bit activations instead of floating points, etc. It was shown, that this approach can lead to a DNN representation, which accuracy is not very far off from the state-of-the-art results.
- Network binarization - this is the extreme case of the parameter quantization technique. Due to a new learning algorithm, Expected Back Propagation, which bases on inferring network with binary weights and neurons through a variational Bayesian approach. Techniques using this approach mainly use the real-valued weights as a reference for their binarization. This idea was further extended to binarize both weights and activations, which was implemented in networks such as BinaryNet and XNOR-Net.

1.3. Goal of the project

This project was created during "Team programming project" classes in 2016/17 and was inspired by Intel.

The main goals were:

- to analyze the results of inference, depending on 3 properties:
 - topology
 - weights
 - feature maps
- to discover way to extend binarized network to achieve accuracy of standard network

Each DNN has its own tolerance of inference precision. There are numerous topologies in neural networks. Each topology indicates the number layers and defines how the neurons are connected. The aim is to study the structures and predict which one will behave the best in our environment. Dealing with 1-bit integer weights is going to strongly affect the complexity of computational algorithms and the size of data. The whole workflow is going to be much faster. 1-bit operations are a way simpler than the 32-bit ones. Furthermore, the size of the inputs will require less memory. That is the perfect solution for less efficient devices.

Unfortunately, those properties will negatively affect on the quality of prediction. Thus, the increase of the depth of chosen net should be also taken into consideration. A feature map is an input for the next level neurons. The loss of data and weights precision can be alleviated by the increased number of the feature maps. We wanted to estimate ratio between number of feature maps in basic DNN and number of feature maps in low arithmetic precision DNN that maintain similar quality of the prediction. To sum up, the research may bear out, that BDNNs have some great properties, which should be investigated much deeper. It may also find out, that they achieve better results in some cases and the present solutions should be replaced with BDNN.

Chapter 2

Architecture overview

2.1. Used framework

We decided to work with one of open source framework that can be used for neural networks. The structure of BNN proposed in Binary Connect [3] and XORNET [2] is very similar to standard convolutional network so there is no need to implement BNN from scratch.

All frameworks that we analyse (TensorFlow, Torch, Caffe) have similar functionalities, for example they have already implemented pooling and affine layers so we can easily reuse them.

From available frameworks we decided to choose TensorFlow [6]. It is an open source software library for numerical computation using data flow graphs. Nodes in the graph represent mathematical operations, while the graph edges represent the multidimensional data arrays (called tensors) transported between them. In case on of neural networks, nodes will represent layers of network. The advantage of TensorFlow over other frameworks is really good community support and solid documentation.

2.2. Our implementation

TODO:Opis implementacji w pythonie

TODO:Opis pythonowego opa

TODO:Opis implementacji w C++

2.3. Used networks and datasets

2.3.1. LeNet on MNIST

LeNet is a small 7-level network invented by Yann LeCun et al. It is one of the most known neural network due to it's simplicity. Almost everyone, who is new to machine learning, begins with implementing it on MNIST dataset. However, it does not mean that LeNet is not a powerful tool. According to Wikipedia, couple of banks are applying LeNet to recognise written digits on cheques. Provided implementation was highly affected by Tensorflow's tutorial.

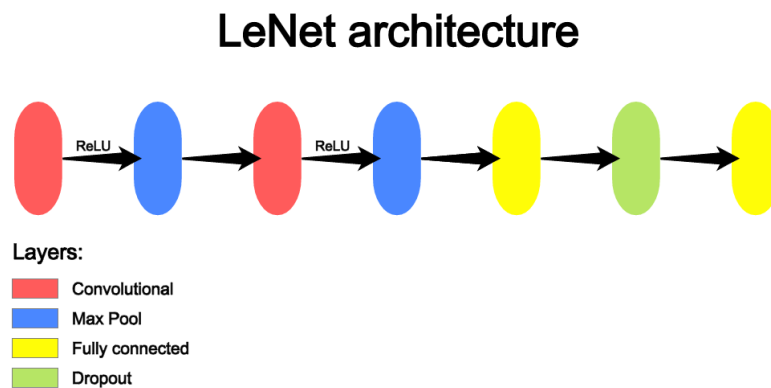
MNIST

MNIST is a subset of NIST. Each example represents a handwritten digit drawn in black&white. Examples were resized and cropped, so that each digit is centered and the size is equal to 18×18 pixels. The dataset is split into 2 parts: training (60,000 pictures) and test (10,000 pictures).

Network architecture

- Convolution layer: kernel size $5 \times 5 \times 32$
- Max Pool layer: kernel size 2×2
- Convolution layer: kernel size $5 \times 5 \times 64$
- Max Pool layer: kernel size 2×2
- Fully-connected layer: 1024 neurons
- Dropout layer: 0.5
- Fully-connected layer: 10 neurons

Figure 2.1: Scheme of LeNet



Parameters

- Weight initialization: truncated normal, stddev 0.1
- Optimizer: Adam Optimizer, learning rate 10^{-4}
- Error: softmax cross entropy with logits
- Non-linearity: ReLU

2.3.2. AlexNet on ImageNet

AlexNet caused a huge breakthrough in deep neural networks after a great success during ImageNet Large-Scale Visual Recognition Challenge in 2012. It achieved top 5 test error rate of 15.4% leaving other opponents far behind. In comparison, the next team achieved an error rate of 26.2%.

2.3.3. Finetuned AlexNet on oxford-102

It has been 5 years since the astonishing victory of AlexNet, although people are still experimenting with it achieving great results. jim-goo(nick z gita, trzeba bedzie zacytowac) proposed a very clever approach of learning, which gave him a top 1 error rate of 7% on oxford-102 set. In this experiment we try to replicate the results and find out whether it is possible to achieve similar accuracy using binarized network.

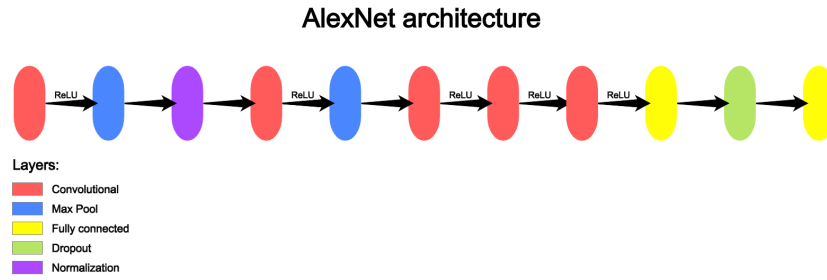
Oxford-flowers-102

Oxford flowers dataset is containing the most common flowers found in United Kingdom. Each image represents a single flower. Within each class, the examples do not differ much, but there exists many similar classes, which makes the task harder. The images are saved in large scale and its' sizes are not normalized. The dataset is split into 3 parts: test (6,149 pictures), train (1,020 images) and validation (1,020 images). In this experiment, the machine is trained on the test set and tested against the train set as the authors of CaffeNet fine-tuned on the Oxford 102 category flower dataset did.

Network architecture

- Convolution layer: kernel size $11 \times 11 \times 96$, learning rate 10^{-3}
- Max Pool layer: kernel size 3×3
- Normalization layer
- Convolution layer: kernel size $5 \times 5 \times 256$, learning rate 10^{-3}
- Max Pool layer: kernel size 3×3
- Convolution layer: kernel size $3 \times 3 \times 384$, learning rate 10^{-3}
- Convolution layer: kernel size $3 \times 3 \times 384$, learning rate 10^{-3}
- Convolution layer: kernel size $3 \times 3 \times 256$, learning rate 10^{-3}
- Fully-connected layer: 4096 neurons, learning rate 10^{-3}
- Dropout layer: 0.5
- Fully-connected layer: 102 neurons, learning rate 10^{-2}

Figure 2.2: Scheme of AlexNet



Parameters

- Weight initialization: pretrained on ImageNet 2012 dataset
- Optimizer: stochastic gradient descent
- Error: softmax cross entropy with logits
- Non-linearity: ReLU

2.3.4. Residual network on CIFAR-10

CIFAR-10

CIFAR-10 is one of the most popular dataset. It is used for object recognition. It includes classes such as: airplanes, automobiles, birds, cats, deers, dogs, frogs, horses, ships and trucks. This dataset contains 60,000 images of size 32×32 . Each class is represented by 6,000 images. The dataset was created by Alex Krizhevsky, Vinod Nair and Geoffrey Hinton. There also exists a bigger dataset called CIFAR-100.

TODO:obrazek, opis sieci, parametry

2.4. Implemented binarization algorithm

2.4.1. Binarized convolution filters from XORNET [2]

In this approach we use special binarized filter for forward propagation.

For an original filter W forward propagation is:

1. Let n be the number of elements in each filter, e.g. if filter is matrix 3×3 , $n = 9$.
2. Let W' be matrix containing signs of W .
3. Calculate A - average of elements for each filter in W' , so A is a vector of size equal to number of filters in layer.
4. Compute standard convolution using matrix W' .
5. Return result of convolution multiplied by A .

The back propagation is standard back-propagation made by using original W matrix. To achieve this result we override standard gradient to change sign to identity.

In this type of binarizations, the gradients are in full precision, therefore the backward-pass still requires convolution between 1-bit numbers and 32-bit floating-points.

2.4.2. Binarized convolution filters inspired by DoReFa [4]

In this approach we do not implement full methods from DoReFa network. We worked on one idea to average filters over all maps at the same time. It is very similar to XORNET scheme.

TODO - opis dokładniejszy co analizuje DoReFa

For original filter W forward propagation is:

1. Let n be the number of elements in each filter, e.g. if filter is matrix 3×3 , $n = 9$.
2. Let W' be matrix containing signs of W .
3. Calculate A average of elements for all filters in W , so A is a scalar.
4. Compute standard convolution using matrix W' .
5. Return result of convolution multiplied by A .

The only difference is that A is now scalar, not a vector. This approach allows to speed up computation for both forward and back propagation (multiplication by scalar is very efficient) comparing to XORNET.

2.4.3. Binarized filters and activations

For original filter W and input I forward propagation is:

1. Let n be the number of elements in each filter, e.g. if filter is matrix 3×3 , $n=9$.
2. Let W_{sign} be matrix containing signs of W divided by A .
3. Calculate A average of elements for each filter in W_{sign} , so A is a vector of size equal to numbers of filters in layer.
4. Let I_{abs} be matrix with absolute value of input.
5. Let I_{sign} be matrix with signs of input.
6. Let K to be result of computation of standard convolution using as input I_{abs} and as weights matrix containing $\frac{1}{n}$ on each position.
7. Compute standard convolution using input I_{sign} and weights W_{sign} .
8. Return result multiplied by K and A .

Chapter 3

Experiments results

3.1. Lenet on MNIST

TODO:wyniki

3.2. AlexNet on Flowers

TODO:wyniki

3.3. Residual network on CIFAR-10

3.3.1. Results and discussion

Binary Resnet XORNET style

Results The results of experiments are in table

TODO:nauczyć się numerować tabele

Table 3.1: Results of Resnet and BinResnet

blocks	layers	epochs	learning rate	ResNet	BinResNet
2	14	0.1	15	80%	77%
5	32	0.1	15	82%	81%
18	110	0.1	20	84%	82%
2	14	0.01	15	78%	76%
5	32	0.01	15	79%	78%
18	110	0.01	20	81%	80%

Discussion The results of binarized network are very good comparing to standard ResNet. The property of ResNet is preserved and there are no much difference in accuracy. That is a very interesting result, because it shows that binarizing ResNet is only slightly decreasing accuracy so it is worth to use it to save memory and computation time, especially on CPU.

Binary Resnet DoReFa style

Results The results of experiments are included in table below.

TODO:nauczyć się ładnie formatować tabele w latexie

Table 3.2: Results of Resnet and BinResnet

blocks	layers	epochs	learning rate	ResNet	BinResNet	BinResnet more epochs
2	14	0.1	15	80%	77%	78% (25 epochs)
5	32	0.1	15	82%	75%	8% (40 epochs)
18	110	0.1	20	84%	79%	80% (40 epochs)
2	14	0.01	15	78%	61%	65%(25 epochs)
5	32	0.01	15	79%	63%	67% (25 epochs)
18	110	0.01	20	81%	62%	67% (40 epochs)

Discussion The network with binarization which takes average over all filters is affected more easily by changes in model parameters. For smaller learning rates it learns quite slow, achieving around 75% of accuracy compared to classical resnet in same number of epochs. It has a potential for longer training - training it for more than 40 epochs can generate accuracy around 0.8.

Of course it is faster than XORNET implementation, but because it requires more training and gets lower accuracy, one must decide if this is good approach based on available computation machine. For testing on personal laptop, XORNET is better solution.

Binary Resnet with binary weights and activation

Results

TODO:zdebugować kod i puścić eksperymenty jeszcze raz

Discussion

- Implementation which is using `tf.nn.conv2d` 2 times in each binarization process is much slower.
- Network with binarized both weights and activation is much harder to learn. In 2-layer Lenet the results were still very stable, but adding more layers, for example 14 in 2-blocks ResNet made network able to learn only during first few rounds of computation.
- Smaller learning rate allows network to learn. Unfortunately the problem of deeper network still occurs, so the advantage of ResNet is not preserved.

Chapter 4

Conclusions

4.1. Discussion of experiments results

TODO:tu opis wspólnych konkluzji dla wszystkich eksperymentów

4.2. Future work

TODO:tu opis co można zrobić jak ma się zmienną 3 bitową, co warto zbadać itd.

Chapter 5

Team members contribution

5.1. Krystyna's contribution

- organising the team's work
- contact with project supervisor
- implementation of Resnet
- experiments on Resnet

5.2. Jakub's contribution

- implementation of TensorFlow binarized convolution
- experiments on Lenet
- implementing C++ gradient for new operation

5.3. Przemysław's contribution

- contact with our partner company
- implementation of TensorFlow binarized convolution
- implementing quality tests for project

5.4. Adam's contribution

- implementation of Alexnet
- experiments on Alexnet
- code review for colleagues

Bibliography

- [1] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, Oriol Vinyals, *Understanding deep learning requires rethinking generalization*, <https://arxiv.org/abs/1611.03530> (2016)
- [2] Mohammad Rastegari, Vicente Ordonez, Joseph Redmon, Ali Farhadi, *XNOR-Net: ImageNet Classification Using Binary Convolutional Neural Networks*, <https://arxiv.org/abs/1603.05279> (2016)
- [3] Matthieu Courbariaux, Yoshua Bengio, Jean-Pierre David, *BinaryConnect: Training Deep Neural Networks with binary weights during propagations*, <https://arxiv.org/abs/1511.00363> (2015)
- [4] Shuchang Zhou, Yuxin Wu, Zekun Ni, Xinyu Zhou, He Wen, Yuheng Zou, *DoReFa-Net: Training Low Bitwidth Convolutional Neural Networks with Low Bitwidth Gradients*, <https://arxiv.org/abs/1606.06160> (2016)
- [5] Kaiming He, Xiangyu Zhang, Shaoqing Ren, Jian Sun, *Deep Residual Learning for Image Recognition*, <https://arxiv.org/abs/1512.03385> (2015)
- [6] <https://www.tensorflow.org/>