



Aalto-yliopisto  
Perustieteiden  
korkeakoulu

# Lossless Compression of Deep Neural Networks (topic-presentation)

*Vilhelm Toivonen*

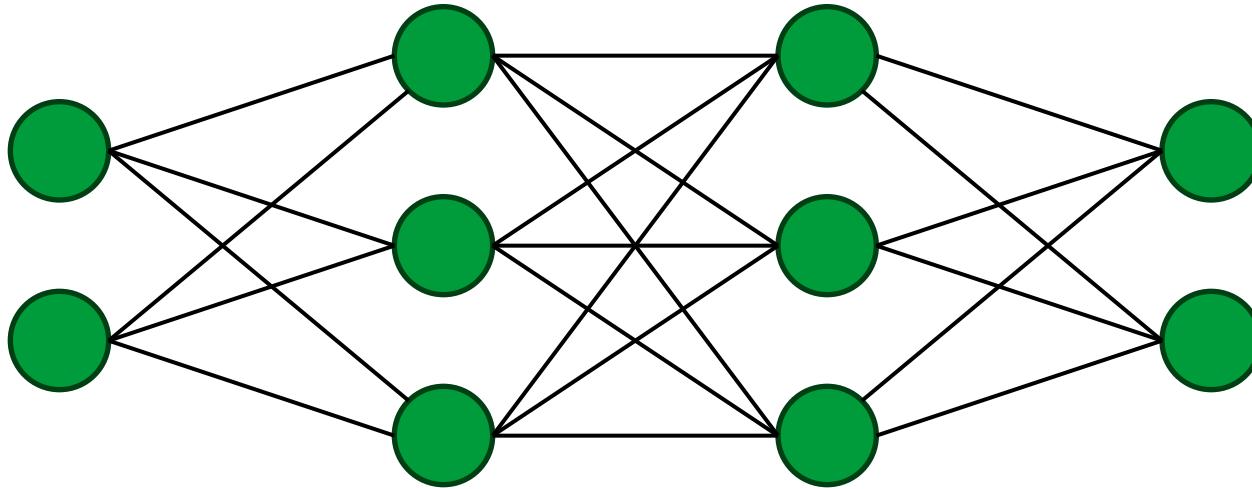
*01.11.2023*

Instructor: *Nikita Belyak*

Supervisor: *Fabricio Oliveira*

Työn saa tallentaa ja julkistaa Aalto-yliopiston avoimilla verkkosivuilla. Muilta osin kaikki oikeudet pidätetään.

# Background – DNNs



Larger DNNs make calculations more intensive

- Slow forward passes
- Computationally expensive to create mathematical programming problems

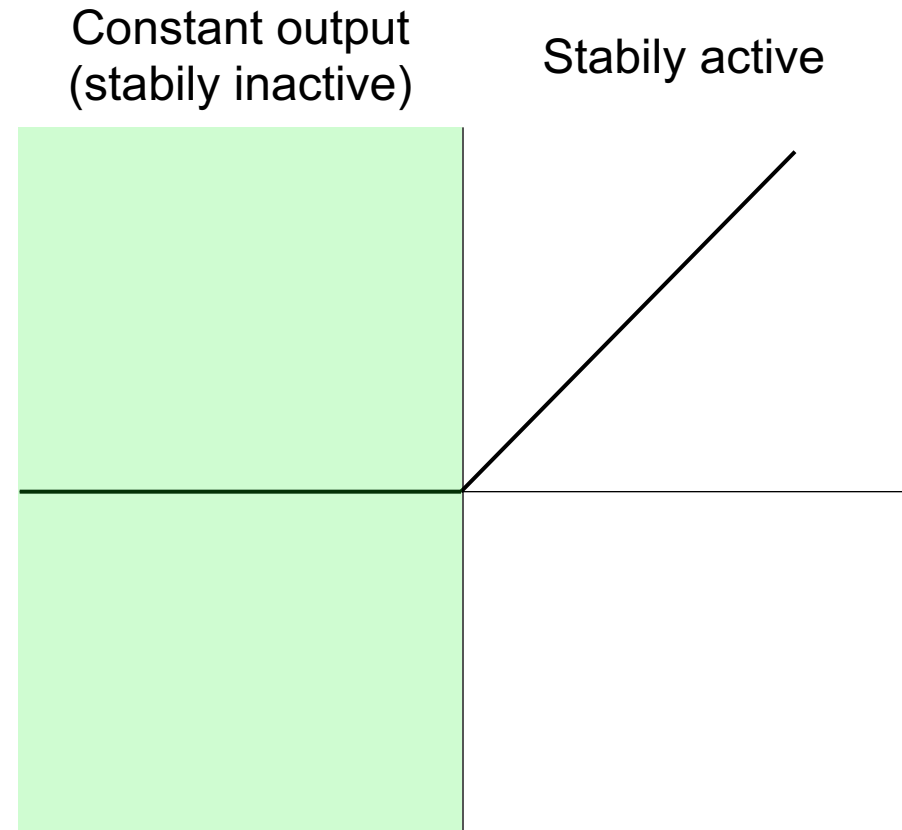
# Background – ReLU

Using the Rectified Linear Unit as the activation, we can separate the output in two parts:

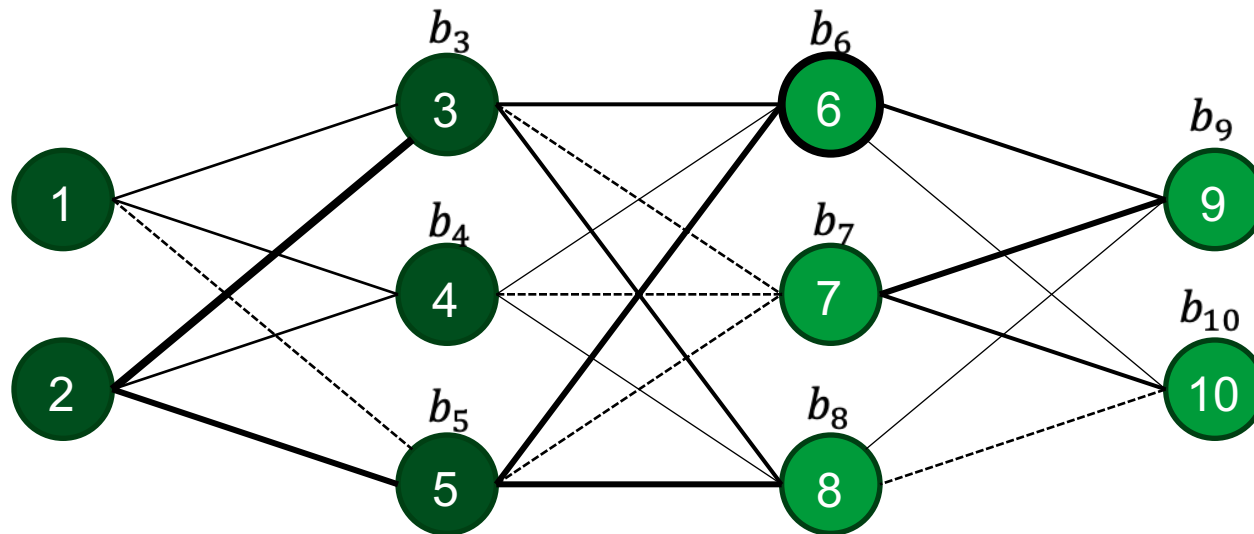
- Constant output
- Stably active

Neurons with constant outputs can be pruned

Stably active neurons can be simplified to system of linear equations, and linearly dependent neurons pruned

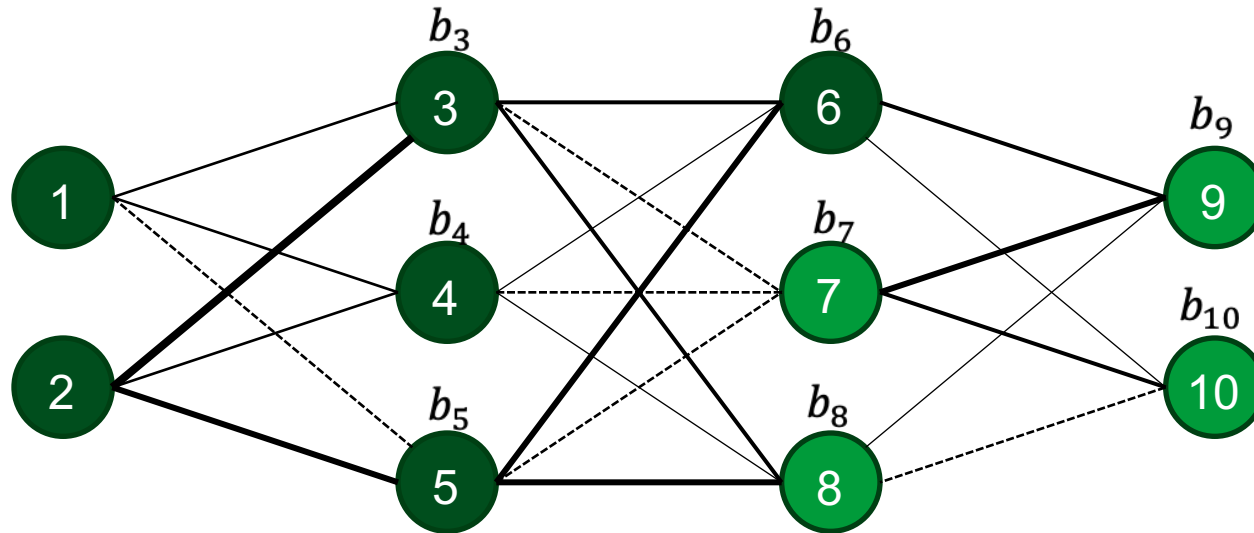


# Background – Example of pruning



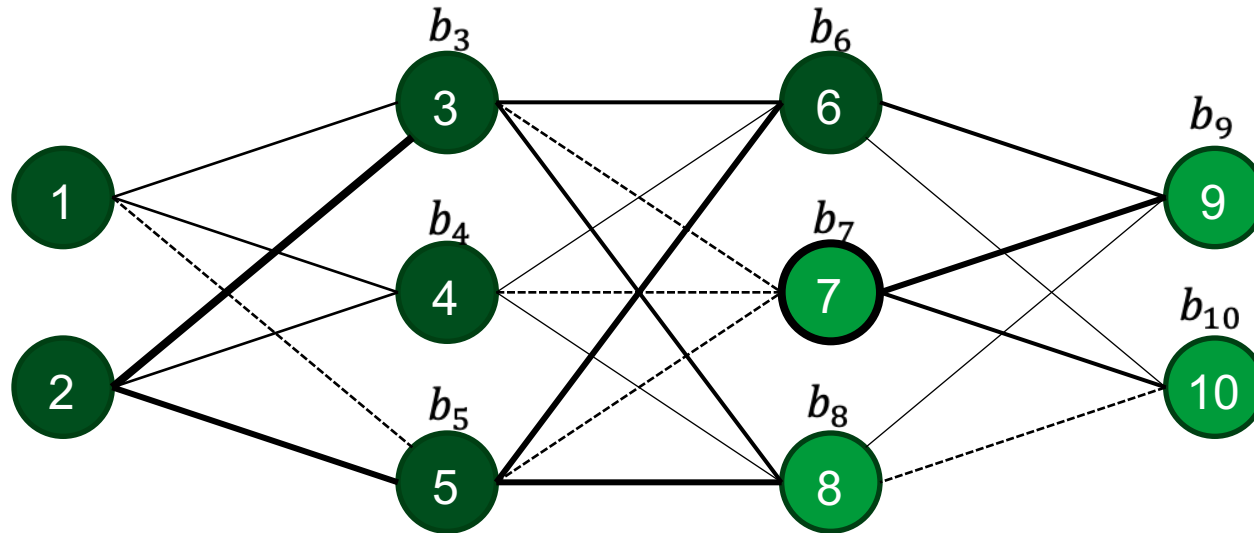
$$\begin{aligned}l &= 2 \\i &= 1 \\ \bar{G}_i^l &> 0 \\ S &= \{\}\end{aligned}$$

# Background – Example of pruning



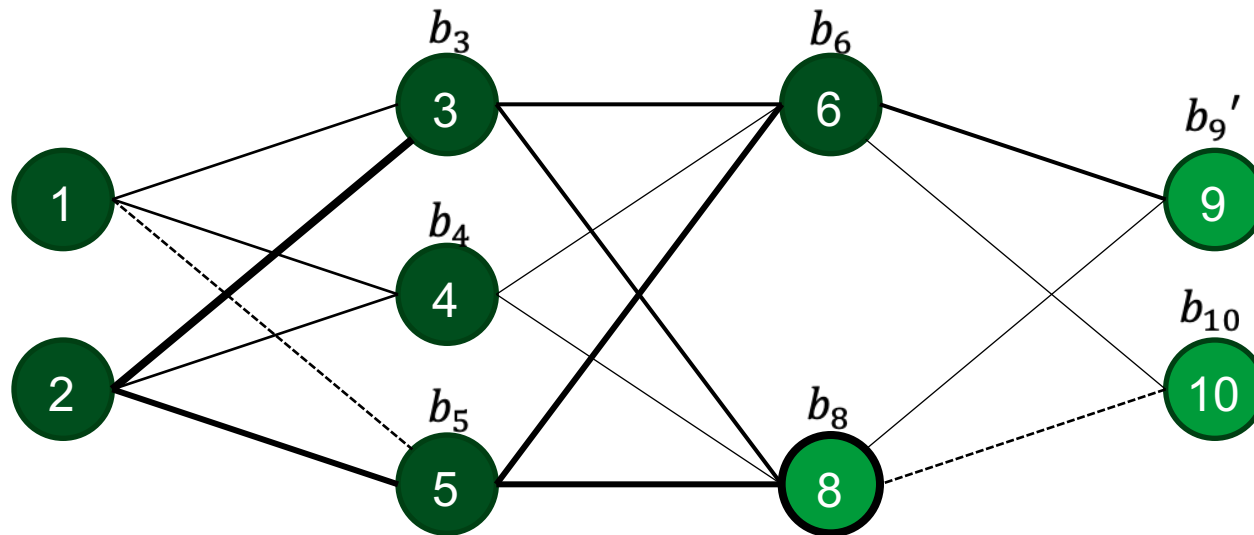
$$\begin{aligned}l &= 2 \\i &= 1 \\ \bar{G}_i^l &> 0 \\ S &= \{1\}\end{aligned}$$

# Background – Example of pruning



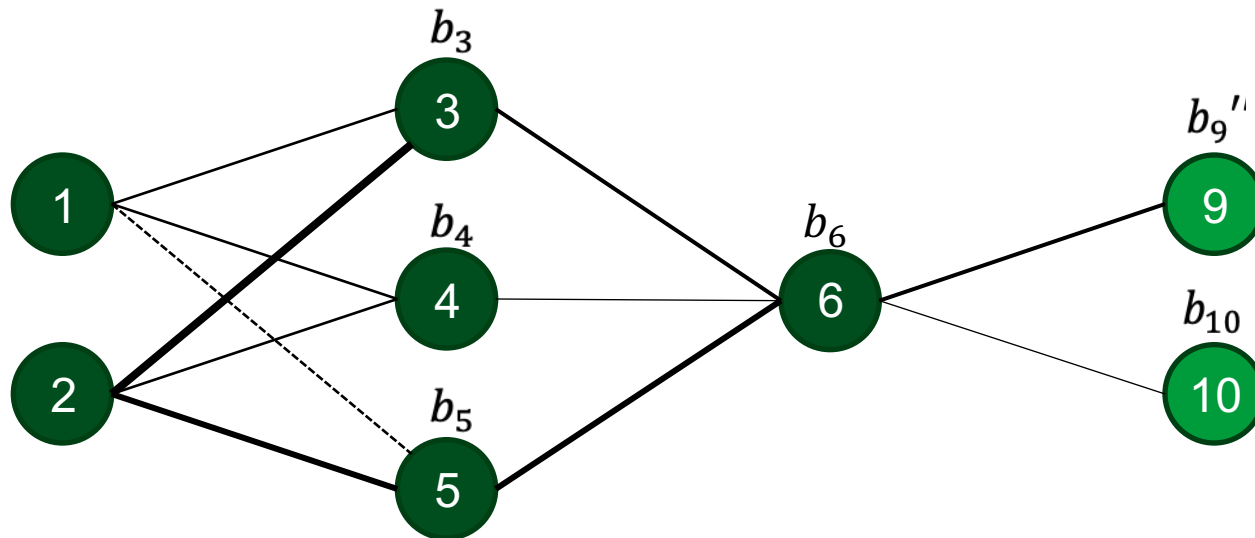
$$\begin{aligned}l &= 2 \\i &= 2 \\W_i^l &= 0 \\S &= \{1\}\end{aligned}$$

# Background – Example of pruning



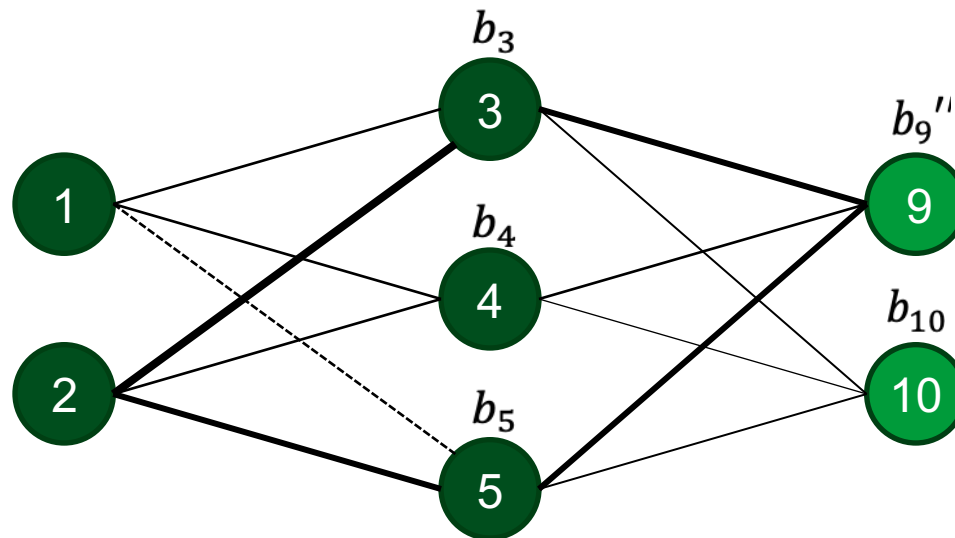
$$l = 2$$
$$i = 3$$
$$\bar{G}_i^l > 0$$
$$S = \{1\}$$

# Background – Example of pruning





# Background – Example of pruning



# Goals

Develop a working algorithm to compress neural networks.

Determine the optimal set of hyperparameters that result in the most effective compression. Consider L1 and L2 norm regularizations with different weights, and at least the SGD, Momentum and Adam optimizers.

```
1: for  $l \leftarrow 1, \dots, L$  do
2:    $S \leftarrow \{\}$  ▷ Set of stable units left in layer  $l$ 
3:   Unstable  $\leftarrow$  False ▷ If there are unstable units in layer  $l$ 
4:   for  $i \leftarrow 1, \dots, n_l$  do
5:     if  $G_i^l < 0$  for  $x \in \mathbb{D}$  or  $W_i^l = 0$  then ▷ Stably inactive, constant output
6:       if  $i < n_l$  or  $|S| > 0$  or Unstable then
7:         if  $W_i^l = 0$  and  $b_i^l > 0$  then
8:           for  $j \leftarrow 1, \dots, n_{l+1}$  do
9:              $b_j^{l+1} \leftarrow b_j^{l+1} + w_{ji}^{l+1} b_i^l$ 
10:          end for
11:        end if
12:        Remove unit  $i$  from layer  $l$  ▷ Unit  $i$  is not necessary
13:      end if
14:    else if  $G_i^l > 0$  for  $x \in \mathbb{D}$  then ▷ Stably active
15:      if  $\text{rank}(W_{S \cup \{i\}}^l) > |S|$  then
16:         $S \leftarrow S \cup \{i\}$  ▷ Keep unit in the network
17:      else
18:        Find  $\{\alpha_k\}_{k \in S}$  such that  $w_i^l = \sum_{k \in S} \alpha_k w_k^l$ 
19:        for  $j \leftarrow 1, \dots, n_{l+1}$  do
20:          for  $k \in S$  do
21:             $w_{jk}^{l+1} \leftarrow w_{jk}^{l+1} + \alpha_k w_{ji}^{l+1}$ 
22:          end for
23:           $b_j^{l+1} \leftarrow b_j^{l+1} + w_{ji}^{l+1} (b_i^l - \sum_{k \in S} \alpha_k b_k^l)$ 
24:        end for
25:        Remove unit  $i$  from layer  $l$  ▷ Unit  $i$  is no longer necessary
26:      end if
27:    else
28:      Unstable  $\leftarrow$  True
29:    end if
30:  end for
31:  if not Unstable then ▷ All units left in layer  $l$  are stable
32:    if  $|S| > 0$  then ▷ The units left have varying outputs
33:      Create matrix  $\bar{W} \in \mathbb{R}^{n_l \times n_{l+1}}$  and vector  $\bar{b} \in \mathbb{R}^{n_{l+1}}$ 
34:      for  $i \leftarrow 1, \dots, n_{l+1}$  do
35:         $\bar{b}_i \leftarrow b_i^{l+1} + \sum_{k \in S} w_{ik}^{l+1} b_k^l$ 
36:        for  $j \leftarrow 1, \dots, n_{l-1}$  do
37:           $\bar{w}_{ij} \leftarrow \sum_{k \in S} w_{kj}^l w_{ik}^{l+1}$ 
38:        end for
39:      end for
40:      Remove layer  $l$ ; replace parameters in next layer with  $\bar{W}$  and  $\bar{b}$ 
41:    else ▷ Only unit left in layer  $l$  has constant output
42:      Compute output  $\Upsilon$  for any input  $\chi \in \mathbb{D}$ 
43:       $(W^{L+1}, b^{L+1}) \leftarrow (0, \Upsilon)$  ▷ Set constant values in output layer
44:      Remove layers 1 to  $L$  and break ▷ Remove all hidden layers and
45:    end if
46:  end if
47: end for
```

# Limitations

Only consider DNNs with the ReLU activation.

Do not consider other architectures, such as CNNs or RNNs.

Do not focus on calculating the bounds for the neurons.

# Sources and materials

- Thiago Serra, Abhinav Kumar, and Srikumar Ramalingam 2020. Lossless Compression of Deep Neural Networks. Bucknell University and the University of Utah. Usa
- Linkola, J. 2023. Reformulating deep neural networks as mathematical programming problems. Bachelor thesis. Aalto-University. School of Science. Espoo.
- ML\_as\_MO package ([https://github.com/gamma-opt/ML as MO](https://github.com/gamma-opt/ML_as_MO)). Includes implementations for calculating bounds

# Tools

**Julia**, includes an implementation for calculating the bounds for neurons.

# Schedule

- Introduction to topic and understanding given materials  
06-10/2023
- Experimentation 10-12/2023
- Topic presentation 11/2023
- Writing the thesis 11/2023-01/2024
- Results presentation 12/2023
- Thesis ready 01/2024