





SPEEDING-UP PRUNING FOR ARTIFICIAL NEURAL NETWORKS: PRESENTING ACCELERATED ITERATIVE MAGNITUDE PRUNING

Marco Zullich Eric Medvet Felice Andrea Pellegrino University of Trieste (Italy) Alessio Ansuini Area Research & Technology (Italy)

# ITERATIVE MAGNITUDE PRUNING: STATE OF THE ART

## PRUNING IN ANN & MAGNITUDE PRUNING

- Pruning a Neural Network → removing parameters from it
- Large number of criteria for pruning
- Magnitude pruning deletes parameters having small magnitude



## PRUNING AND RE-TRAINING

- A simple application of pruning degrades the ANN performance
- After pruning, a **re-training phase** follows
- Re-training is operated only on parameters having survived the pruning



**ITERATIVE MAGNITUDE PRUNING** (IMP)

## MAIN TECHNIQUES FOR RE-TRAINING



## PROSAND CONS OF WR & LRR

#### PROS

Reach very high pruning rates (>95%) with **similar or better performance** w.r.t. unpruned network

#### CONS

Especially if compared to other methods, requires application of many sequential iterations

If a target sparsity is known from the beginning, is it possible to fastforward the execution of IMP for all the iterations but the last one?

# PRESENTING ACCELERATED ITERATIVE MAGNITUDE PRUNING

## ACCELERATING IMP

- Unpruned ANN trained for T epochs
- Prune for K iterations
- Iterations 1, ..., K 1: retrain for  $\tau$  epochs,  $\tau \ll T$
- Accelerated Iterative
  Magnitude Pruning (AIMP)
- Test with VGG-19 on CIFAR10 dataset
- *T* = 160; *K* = 20; *p* = 0.2



## DRAWBACKS & DIRECTIONS FOR FUTURE WORK

- Trials on IMP + LRR were not as satisfying as IMP + WR
- Median accuracy: 93.68 % VS. 63.62 % (τ = 50)

No proper criterion to determine an optimal  $\tau$ 

• AIMP seems to work only when **overall pruning rate** is very high ( $\geq$  98%)

ICPR 2020 – Zullich, Pellegrino, Medvet, Ansuini



## Thanks for the attention!

# Contacts: marco.zullich@phd.units.it



