

Mixed-Precision Algorithms for Training Deep Neural Networks

Lars Ruthotto

Abstract

Abstract. We present ongoing research on improving the efficiency of deep neural network training through mixed-precision computation and modified Gauss-Newton algorithms. By adjusting the precision of various computations and parameters, we achieve reduced model sizes and lower computational and communication costs. To minimize the expense of Gauss-Newton approximations, we utilize advanced automatic differentiation techniques. We demonstrate our results in various learning tasks, including physics-informed neural networks, supervised learning, and generative modeling through normalizing flows and flow matching.