

Initial Guessing Bias: Connecting Initial Predictive Behavior to Trainability in Deep Neural Networks

Abstract

The initial state of deep neural networks (DNNs) significantly influences the subsequent learning process. In classification problems, DNNs can exhibit an intrinsic predictive bias toward certain classes before any training occurs—a phenomenon known as Initial Guessing Bias (IGB). We present a theoretical analysis revealing how architectural design choices—including activation functions, pooling layers, normalization schemes and network depth—shape IGB.

Furthermore, we establish a rigorous connection between the predictive behaviors described by the IGB framework and Mean-Field (MF) theories, which characterize the trainability of deep networks via order-to-chaos phase diagrams.

This correspondence demonstrates that the initial predictive behaviour directly relates to the conditions required for stable and efficient learning. Counter-intuitively, optimal initialization for training does not correspond to an unbiased or neutral initial state but rather to a state with transient yet substantial predictive bias. This analysis offers a unified framework that deepens our theoretical understanding of how the initial state of DNNs governs both predictive bias and trainability.

Speaker

Emanuele Francazi

PhD student

École Polytechnique Fédérale de
Lausanne

