Using Fast Weights to Attend to the Recent Past

Improving neural activity representations and weights to capture regularities among inputs restricts development of artificial neural systems. A natural alternative to prevent this resitrction from arising is leveraging variable time-scale operations motivated by dynamics of synapses. Weights that change faster than standard weights and slower than neural activities present a neurally plausible scheme for attending to events in the short term. The work build on this insight and makes use of fast weights as an attention mechanism to learn neural activity patterns while preventing the need to store their large copies. Fast weight matrices constructed using hidden representations are leveraged efficiently without computing the complete weight matrix and allow learning models to efficiently attend towards complex inputs in different learning settings.

Utilization of associative memory is an essential component in scenarios wherein the model may need to attent to large memory inputs. These inputs can be represented as compact synaptic footprints which can be learnt efficiently using a faster temporal scheme. Evidence of fast neural computation may be observed in physiolog wherein the human brain leverages a variety of short-term plasticity mechanisms that operate on different timescales. This gives rise to fast weight updates in artificial neural networks wherein fast weight matrices are updated at an increased timescale in comparison to standard weights. Fast weights update recent history using weighted outer products of hidden representations which serve as attention weights to the model. In order to eliminate extra computation of fast weights, outer products are directly computed using previous hidden representations. Stabilization of the product is carried out using Layer Normalization which prevents the hidden matrices from vanishing or exploding values.