

Long-Term Dependencies: The Challenge

Sargur Srihari
srihari@buffalo.edu

Topics in Sequence Modeling

- Recurrent Neural Networks
 1. Unfolding Computational Graphs
 2. Recurrent Neural Networks
 3. Bidirectional RNNs
 4. Encoder-Decoder Sequence-to-Sequence Architectures
 5. Deep Recurrent Networks
 6. Recursive Neural Networks
 7. The Challenge of Long-Term Dependencies
 8. Echo-State Networks
 9. Leaky Units and Other Strategies for Multiple Time Scales
 10. LSTM and Other Gated RNNs
 11. Optimization for Long-Term Dependencies
 12. Explicit Memory

Challenge of Long-Term Dependencies

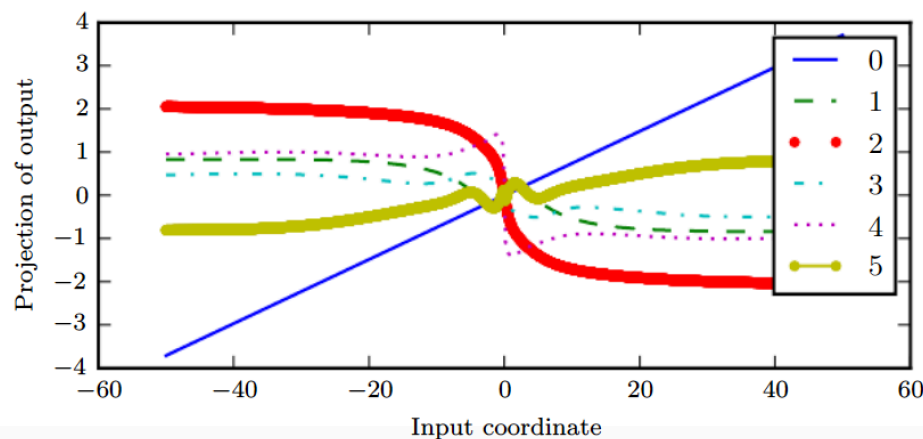
- Neural network optimization face a difficulty when computational graphs become deep, e.g.,
 - Feedforward networks with many layers
 - RNNs that repeatedly apply the same operation at each time step of a long temporal sequence
- Gradients propagated over many stages tend to either vanish (most of the time) or explode (damaging optimization)
- The difficulty with long-term dependencies arise from exponentially smaller weights given to long-term interactions (involving multiplication of many Jacobians)

Vanishing and Exploding Gradient Problem

- Suppose a computational graph consists of repeatedly multiplying by a matrix W
- After t steps this is equivalent to multiplying by W^t
- Suppose W has an eigendecomposition $W = V \text{diag}(\lambda) V^{-1}$
 - In this case it is straightforward to see that
$$W^t = (V \text{diag}(\lambda) V^{-1})^t = V \text{diag}(\lambda)^t V^{-1}$$
 - Any eigenvalues λ_i that are not near an absolute value of 1 will either explode if they are greater than 1 in magnitude and vanish if they are less than 1 in magnitude
- Vanishing gradients make it difficult to know which direction the parameters should move to improve cost
- Exploding gradients make learning unstable

Function Composition in RNNs

- RNNs involve composition of the same function multiple times, one per step
- These compositions can result in extremely nonlinear behavior



Composing many nonlinear functions:

\tanh here

\mathbf{h} has 100 dimensions
mapped to a single
dimension

Most of the space it has a
small derivative and highly
nonlinear elsewhere

- Problem particular to RNNs

A scalar case and a possible solution

- In the scalar case, imagine multiplying w by itself many times
- The product w^t will either vanish or explode depending on the magnitude of w
- However if we make a nonrecurrent network that has a different weight w^t at each step the situation is different
- If the initial state is 1, the state at time t is $\prod_t w^{(t)}$
- If $w^{(t)}$ values are generated randomly, independently with zero mean and variance v , then the variance of the product is $O(v^n)$
- To obtain some desired variance v^* , we may choose individual weights with variance $\sqrt[n]{v^*}$