

Optimization for Training Deep Models

Jinhwan Suk

Department of Mathematical Science, KAIST

May 28, 2020

Contents

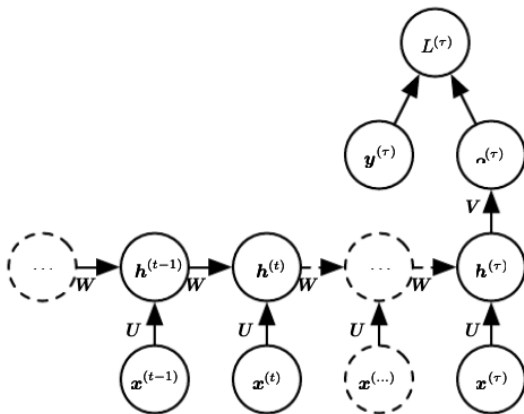
- 1 Recursive Neural Network
- 2 The Challenge of Long-Term Dependencies
- 3 Echo State Networks
- 4 Leaky Units and Other Strategies for Multiple Time Scales
- 5 The Long Short-Term Memory and Other Gated RNNs
- 6 Optimization for Long-Term Dependencies
- 7 Attention model

Contents

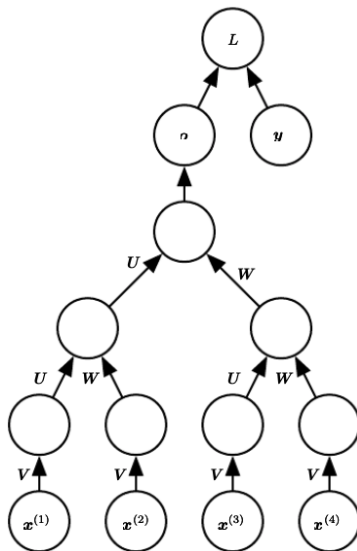
- 1 Recursive Neural Network
- 2 The Challenge of Long-Term Dependencies
- 3 Echo State Networks
- 4 Leaky Units and Other Strategies for Multiple Time Scales
- 5 The Long Short-Term Memory and Other Gated RNNs
- 6 Optimization for Long-Term Dependencies
- 7 Attention model

Recursive Neural Network \neq Recurrent Neural Network (RNN)

Recursive Neural Network



Recursive Neural Network



The depth can be reduced from $O(\tau)$ to $O(\log \tau)$
 \Rightarrow long-term dependency \downarrow

Best structure of the tree??

- ① fixed structure e.g. Balanced binary tree
- ② unfixed structure
 - 1) use natural language parser 2) the learner to discover

Contents

- 1 Recursive Neural Network
- 2 The Challenge of Long-Term Dependencies
- 3 Echo State Networks
- 4 Leaky Units and Other Strategies for Multiple Time Scales
- 5 The Long Short-Term Memory and Other Gated RNNs
- 6 Optimization for Long-Term Dependencies
- 7 Attention model

The Challenge of Long-Term Dependencies

- RNN : $h^{(t)} = \sigma(b + W^T h^{(t-1)} + Ux^{(t)}), t = 1, 2, \dots, \tau$

$$\frac{\partial h^{(t)}}{\partial h^{(t-1)}} = h^{(t)}(1 - h^{(t)})W$$

$$\frac{\partial h^{(t)}}{\partial h^{(t-n)}} = \left(\prod_{i=t-n+1}^t h^{(i)}(1 - h^{(i)}) \right) W^n$$

- If n is large, power of $h^{(t-n)}$ can be extremely big or small.

“The game became interesting as the players warmed up although it was boring for the first half.”

Contents

- 1 Recursive Neural Network
- 2 The Challenge of Long-Term Dependencies
- 3 Echo State Networks**
- 4 Leaky Units and Other Strategies for Multiple Time Scales
- 5 The Long Short-Term Memory and Other Gated RNNs
- 6 Optimization for Long-Term Dependencies
- 7 Attention model

Learning parameters : hidden \rightarrow hidden, hidden \rightarrow output

only learn the output weights \approx kernel machine

Q : How do we initialize the weights so that a rich set of a histories can be represented in the output state?

A : spectral radius > 1

⚠ Consider saturation problem

Contents

- 1 Recursive Neural Network
- 2 The Challenge of Long-Term Dependencies
- 3 Echo State Networks
- 4 Leaky Units and Other Strategies for Multiple Time Scales**
- 5 The Long Short-Term Memory and Other Gated RNNs
- 6 Optimization for Long-Term Dependencies
- 7 Attention model

Leaky Units and Other Strategies for Multiple Time Scales

- Adding Skip Connections through Time
- Removing Connections
- Leaky Units and a Spectrum of Different Time Scales

$$h^{(t)} \leftarrow \alpha h^{(t-1)} + (1 - \alpha) \sigma(W^T h^{(t-1)} + Ux^{(t)} + b)$$

\Rightarrow Accumulate information

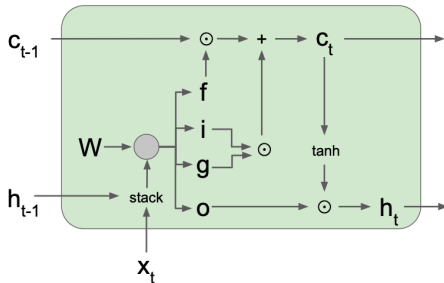
Contents

- 1 Recursive Neural Network
- 2 The Challenge of Long-Term Dependencies
- 3 Echo State Networks
- 4 Leaky Units and Other Strategies for Multiple Time Scales
- 5 The Long Short-Term Memory and Other Gated RNNs
- 6 Optimization for Long-Term Dependencies
- 7 Attention model

The Long Short-Term Memory and Other Gated RNNs

LSTM

- Solve vanishing gradient problem
- Once that information has been used, it might be useful to forget the old state.
- Mechanism to forget the old state?



The Long Short-Term Memory and Other Gated RNNs

LSTM

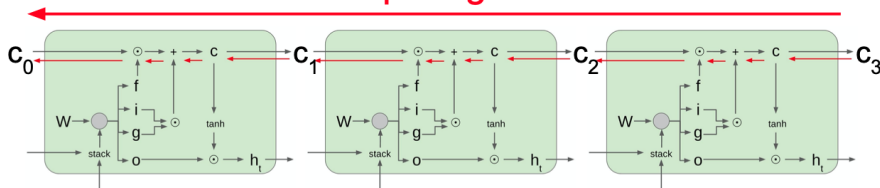
- **Memory cell** c_t is controlled by a forget gate f and input gate g

$$c^{(t)} = f \odot c^{(t-1)} + g \odot \sigma(b + Wh^{(t-1)} + Ux^{(t)})$$

- The output $h^{(t)}$ is controlled by c_t and output gate o

$$h^{(t)} = o \odot \tanh(c^{(t)})$$

Uninterrupted gradient flow!



The Long Short-Term Memory and Other Gated RNNs

Other Gated RNNs

$$h^{(t)} = u_t \odot h^{(t-1)} + (1 - u_t) \odot \sigma \left(b + W(r_t \odot h^{(t-1)}) + Ux^{(t)} \right)$$

- \mathbf{u} stands for “update” gate and \mathbf{r} for “reset” gate.
- $u_t = \sigma \left(b^u + W^u h^{(t-1)} + U^u x^{(t)} \right)$
- $r_t = \sigma \left(b^r + W^r h^{(t-1)} + U^r x^{(t)} \right)$
- faster computation, less parameter

Contents

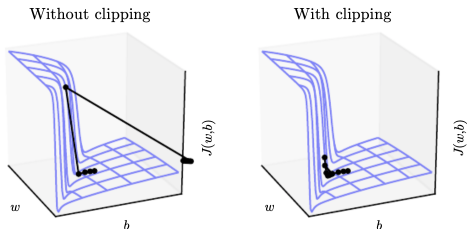
- 1 Recursive Neural Network
- 2 The Challenge of Long-Term Dependencies
- 3 Echo State Networks
- 4 Leaky Units and Other Strategies for Multiple Time Scales
- 5 The Long Short-Term Memory and Other Gated RNNs
- 6 Optimization for Long-Term Dependencies**
- 7 Attention model

Optimization for Long-Term Dependencies

- Vanishing and exploding gradient problems
 - Second-order method, BFGS
 - Nesterov momentum + careful initialization
 - **LSTM + SGD**

Optimization for Long-Term Dependencies

- Clipping Gradients



if $\|g\| > \nu$, then $g \leftarrow \frac{g\nu}{\|g\|}$

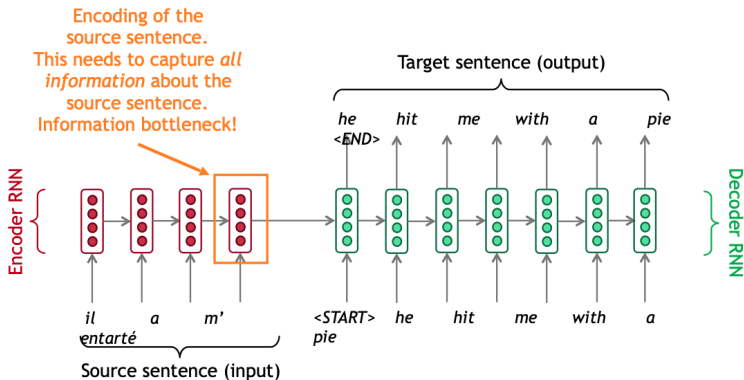
- Regularization term

$$\Omega = \sum_t \left(\frac{\left\| (\nabla_{h^{(t)}} L) \frac{\partial h^{(t)}}{\partial h^{(t-1)}} \right\|}{\left\| \nabla_{h^{(t)}} L \right\|} - 1 \right)^2$$

Contents

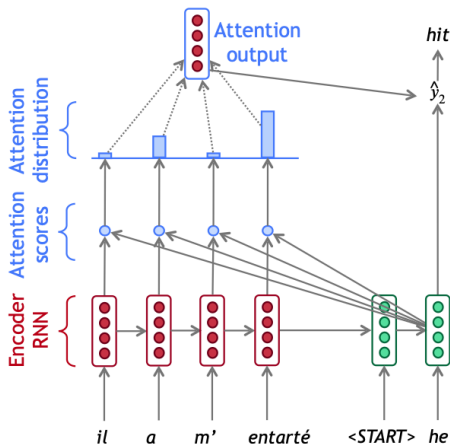
- 1 Recursive Neural Network
- 2 The Challenge of Long-Term Dependencies
- 3 Echo State Networks
- 4 Leaky Units and Other Strategies for Multiple Time Scales
- 5 The Long Short-Term Memory and Other Gated RNNs
- 6 Optimization for Long-Term Dependencies
- 7 Attention model

Attention model



Attention model

Core Idea : on each step of the decoder, use direct connection to the encoder to **focus** on a particular part of the source sequence



Attention model

Attention is all you need (Transformer)

- Non-recurrent sequence model
- Instead, add “positional encodings” to the input embedding
- Self-attention

$$Attention(Q, K, V) = softmax(\frac{QK^T}{\sqrt{d_k}})V$$

- Multi-Head Attention

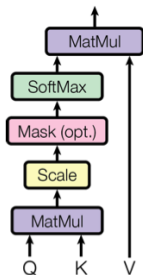
$$MultiHead(Q, K, V) = Concat(head_1, \dots, head_h)W^o$$

$$head_i = Attention(QW_i^Q, KW_i^K, VW_i^V)$$

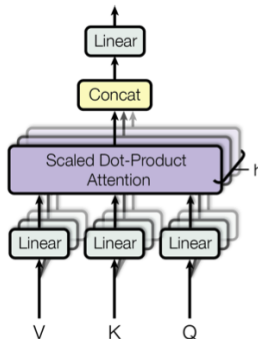
Attention model

Attention is all you need (Transformer)

Scaled Dot-Product Attention



Multi-Head Attention



Attention model

Attention is all you need (Transformer)

