

On the difficulty of training Recurrent Neural Networks

By R. Pascanu, T. Mikolov, Yoshua Bengio

Ludwig Winkler

TU Berlin

February 22, 2017

Outline

Notation

Recurrent Neural Networks

Error Propagation in Recurrent Neural Networks

Analytical Analysis of the Gradients

- Power Iteration Method

- Analytical Analysis

Alternative Interpretations

Dealing with the gradients

- Gradient-based Methods

- LSTM

- Neural Turing Machine

Notation

Input at time t	\mathbf{x}_t
Hidden state at time t	\mathbf{h}_t
Output at time t	\mathbf{y}_t
Input weights	\mathbf{W}_1
Recurrent weights	\mathbf{W}_{rec}
Output weights	\mathbf{W}_2
Biases	$\mathbf{b}_1, \mathbf{b}_2$
Model parameters	$\theta = \{\mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_{rec}, \mathbf{b}_1, \mathbf{b}_2\}$
Activation function	$\sigma(\mathbf{x}_t)$
Cost function at time t	ε_t
'Immediate' derivative at time t	$\frac{\partial^+}{\partial \theta}$

Recurrent Neural Network

- Neural Network with constant connections through time
- Hidden state input \mathbf{x}_t and \mathbf{h}_{t-1}
- Retains past input information in \mathbf{h}_t
- For a simple three layer network the hidden states are updated as

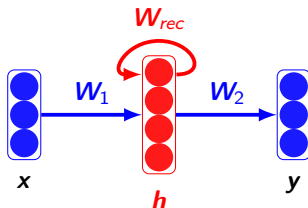
$$\mathbf{h}_t = \mathcal{F}(\mathbf{x}_t, \mathbf{h}_{t-1}; \theta)$$

$$\mathbf{h}_t = \sigma(\mathbf{W}_{rec} \cdot \mathbf{h}_{t-1} + \mathbf{W}_1 \cdot \mathbf{x}_t + \mathbf{b}_1)$$

$$\mathbf{y}_t = \sigma(\mathbf{W}_2 \cdot \mathbf{h}_t)$$

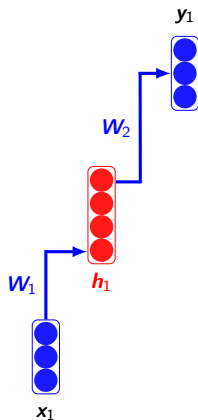
Recurrent Neural Network

- Neural Network with constant connections through time
- Hidden state input \mathbf{x}_t and \mathbf{h}_{t-1}
- Retains past input information in \mathbf{h}_t
- For a simple three layer network the hidden states are updated as



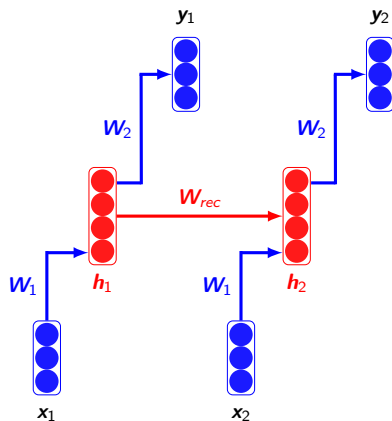
Recurrent Neural Network

- RNNs can be unfolded to visualize their recurrent information flow



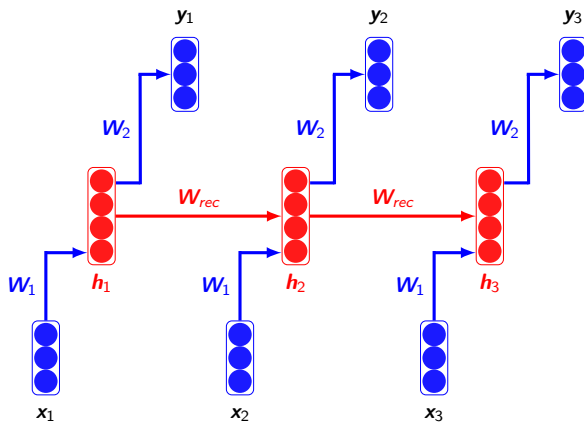
Recurrent Neural Network

- RNNs can be unfolded to visualize their recurrent information flow



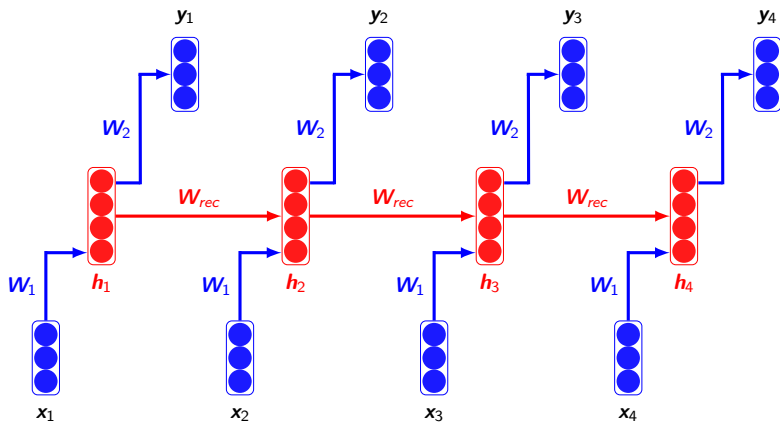
Recurrent Neural Network

- RNNs can be unfolded to visualize their recurrent information flow



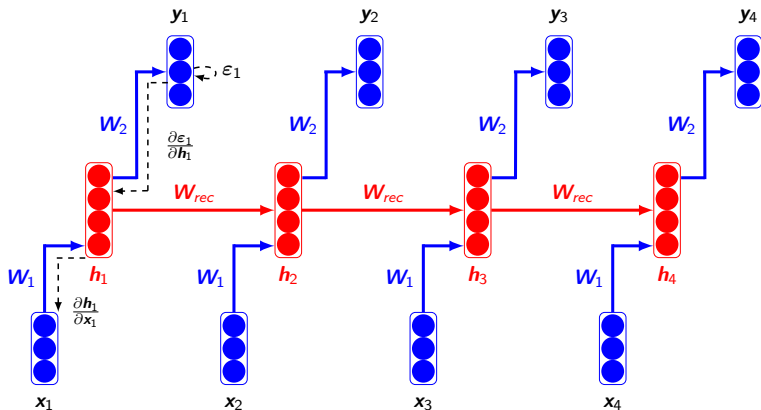
Recurrent Neural Network

- RNNs can be unfolded to visualize their recurrent information flow



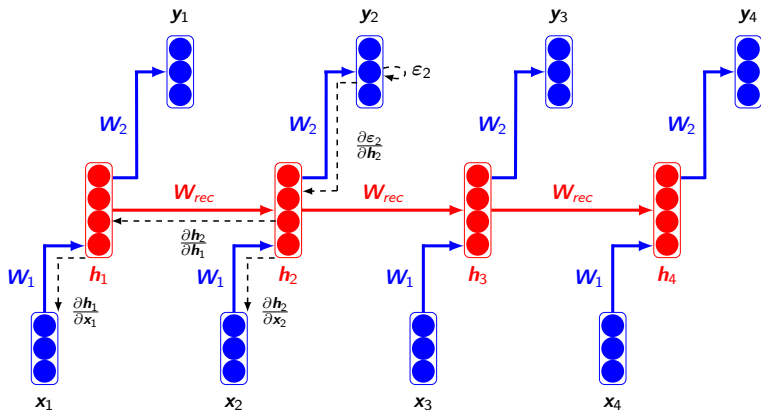
Error Propagation in Unfolded RNN

- Trained with Backpropagation Through Time (BPTT)



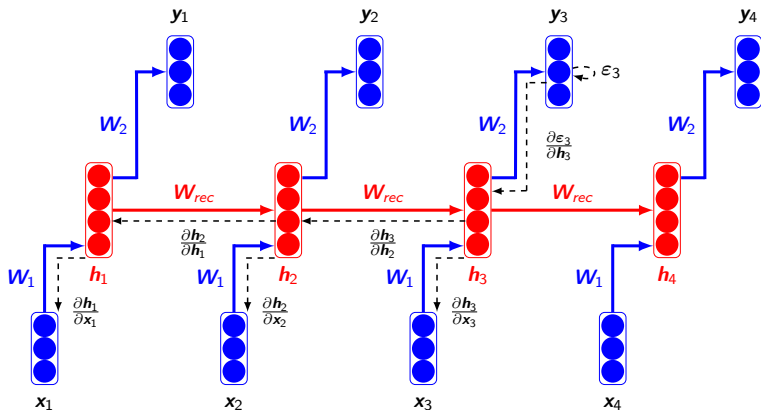
Error Propagation in Unfolded RNN

- Trained with Backpropagation Through Time (BPTT)



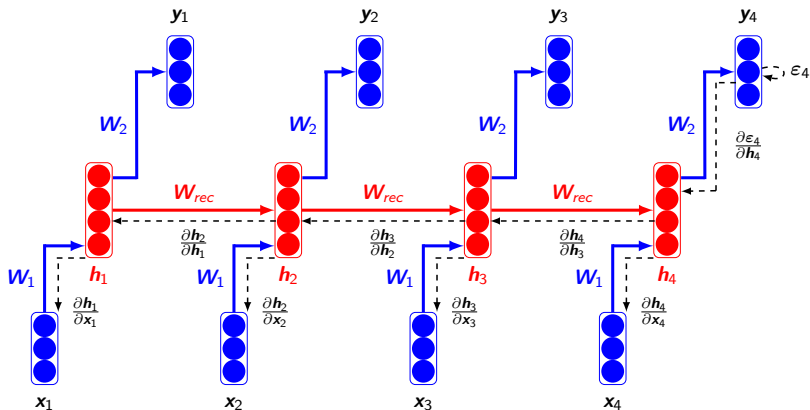
Error Propagation in Unfolded RNN

- Trained with Backpropagation Through Time (BPTT)



Error Propagation in Unfolded RNN

- Trained with Backpropagation Through Time (BPTT)



Error Propagation in RNN

- An unfolded RNN can be interpreted as a deep neural network
- The cost \mathcal{E} can be decomposed into its temporal contributions

$$\frac{\partial \mathcal{E}}{\partial \theta} = \sum_{1 \leq t \leq T} \frac{\partial \mathcal{E}_t}{\partial \theta}$$

- Gradient for all 'temporal' and 'spatial' parameters of the model

Error Propagation in RNN

- Chain rule for the gradients over all past time steps k

$$\frac{\partial \mathcal{E}_t}{\partial \theta} = \sum_{1 \leq k \leq t} \left(\frac{\partial \mathcal{E}_t}{\partial \mathbf{h}_t} \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_k} \frac{\partial^+ \mathbf{h}_k}{\partial \theta} \right)$$

- Cost function for the temporal and spatial parameters of the model
- $\frac{\partial^+ \mathbf{h}_k}{\partial \theta}$ is the 'immediate' derivative disregarding recurrent connections more than one time step away

Error Propagation in RNN

- $\left. \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_k} \right|_{1 \leq k \leq t}$ calculates the gradients for the time steps $k \leq t$

$$\left. \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_k} \right|_{1 \leq k \leq t} = \prod_{k < i \leq t} \frac{\partial \mathbf{h}_i}{\partial \mathbf{h}_{i-1}} = \prod_{k < i \leq t} \mathbf{W}_{rec}^T \text{diag}(\sigma'(\mathbf{h}_{i-1}))$$

- Propagates the error to time step k given the cost at time step t

Analytical Analysis of the Gradients

- Using the identity function for simplification we obtain

$$\mathbf{h}_t = \mathbf{W}_{rec} \cdot \mathbf{h}_{t-1} + \mathbf{W}_1 \mathbf{x}_t + \mathbf{b}_1$$

- The derivative for the $\ell = t - k$ time steps of BPTT is

$$\left. \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_k} \right|_{1 \leq k \leq t} = \prod_{k < i \leq t} \frac{\partial \mathbf{h}_i}{\partial \mathbf{h}_{i-1}} = \left(\mathbf{W}_{rec}^T \right)^\ell$$

Power Iteration Method

- Iterative eigenvector algorithm that will result in \mathbf{v}_1 with $|\lambda_1| > |\lambda_i|$

$$\mathbf{b}_{k+1} = \frac{\mathbf{A}\mathbf{b}_k}{\|\mathbf{A}\mathbf{b}_k\|}$$

- Taking the recurrent relation into account we can reformulate the method as

$$\mathbf{b}_k = \frac{\mathbf{A}^k \mathbf{b}_0}{\|\mathbf{A}^k \mathbf{b}_0\|}$$

Power Iteration Method

- Let λ_i and \mathbf{v}_i be the ordered eigenvalues and their corresponding eigenvectors of matrix \mathbf{A}

$$\mathbf{b}_0 = c_1 \mathbf{v}_1 + c_2 \mathbf{v}_2 + \cdots + c_m \mathbf{v}_m$$

- The product $\mathbf{A}^k \mathbf{b}_0$ can be rewritten with $\mathbf{A}^k \mathbf{v}_i = \lambda_i^k \mathbf{v}_i$ as

$$\begin{aligned} \mathbf{A}^k \mathbf{b}_0 &= c_1 \mathbf{A}^k \mathbf{v}_1 + c_2 \mathbf{A}^k \mathbf{v}_2 + \cdots + c_m \mathbf{A}^k \mathbf{v}_m \\ &= c_1 \lambda_1^k \mathbf{v}_1 + c_2 \lambda_2^k \mathbf{v}_2 + \cdots + c_m \lambda_m^k \mathbf{v}_m \\ &= c_1 \lambda_1^k \mathbf{v}_1 + \lambda_1^k \sum_{j=2}^m c_j \underbrace{\left(\frac{\lambda_j}{\lambda_1} \right)^k}_{|\lambda_j/\lambda_1| < 1} \mathbf{v}_j \approx c_1 \lambda_1^k \mathbf{v}_1 \end{aligned}$$

Analytical Analysis

- Eigendecomposition of \mathbf{W}_{rec}^T gives eigenvalues λ_i and their orthogonal eigenvectors \mathbf{q}_i
- Using \mathbf{q}_i we can linearly decompose $\frac{\partial \mathcal{E}_t}{\partial \mathbf{h}_t}$

$$\begin{aligned}
 \frac{\partial \mathcal{E}_t}{\partial \mathbf{h}_t} &= \sum_{i=1}^N c_i \mathbf{q}_i^T \\
 \sum_{1 \leq k \leq t} \frac{\partial \mathcal{E}_t}{\partial \mathbf{h}_t} \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_k} &= \sum_{i=1}^N c_i \mathbf{q}_i^T \left(\mathbf{W}_{rec}^T \right)^\ell = \sum_{i=1}^N c_i \left(\underbrace{\mathbf{W}_{rec}^\ell \mathbf{q}_i}_{=\lambda_i^\ell \mathbf{q}_i} \right)^T \\
 &= c_1 \lambda_1^\ell \mathbf{q}_1^T + \lambda_1^\ell \sum_{j=2}^N c_j \underbrace{\left(\frac{\lambda_j}{\lambda_1} \right)^\ell}_{|\lambda_j/\lambda_1| < 1} \mathbf{q}_j^T \approx c_1 \lambda_1^\ell \mathbf{q}_1^T
 \end{aligned}$$

Analytical Analysis

Non-Unitarian Case

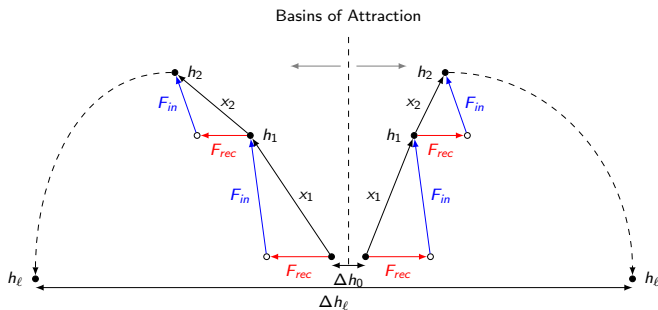
- $\sigma'(\mathbf{h}_k)$ and λ_1 of \mathbf{W}_{rec} are bounded by $\|\text{diag}(\sigma'(\mathbf{h}_k))\| \leq \gamma$
and $\lambda_1 < \frac{1}{\gamma}$

$$\left\| \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} \right\| \leq \|\mathbf{W}_{\text{rec}}\| \|\text{diag}(\sigma'(\mathbf{h}_{t-1}))\| < \frac{1}{\gamma} \gamma < 1$$

- *Sufficient* that $\lambda_1 < 1$ for gradients to vanish
- *Necessary* that $\lambda_1 > 1$ for gradients to explode

Similarities to Dynamical Systems

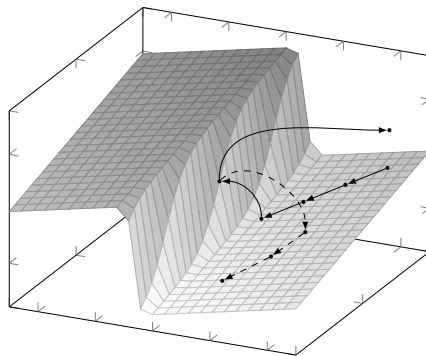
- System approaches attractor asymptotically if in basin of attraction
- Small differences in initialization of \mathbf{W}_{rec} lead to different attractors



$$F_{RNN}(x) = F_{\text{rec}}(h) + F_{\text{in}}(x)$$

Geometric Interpretation

- $\frac{\partial \mathcal{E}_t}{\partial \theta}$ tends to explode along \mathbf{q}_1 of \mathbf{W}_{rec}
- Error surface becomes very steep if \mathbf{q}_1 and $\frac{\partial \mathcal{E}_t}{\partial \theta}$ are aligned
- Regularization needed to prevent excessive gradient descent steps



Gradient-based Methods

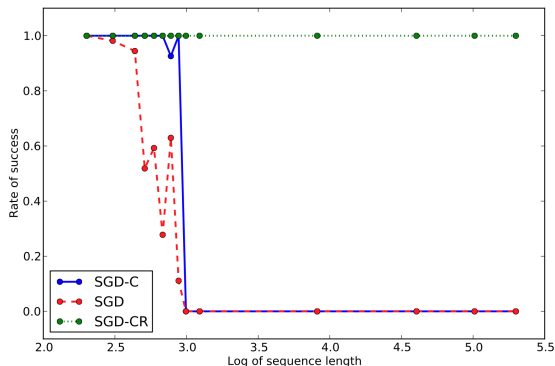
- Clipping exploding gradients if $\|\mathbf{g}\| \geq K$ with $\mathbf{g} \leftarrow K \frac{\mathbf{g}}{\|\mathbf{g}\|}$
- Vanishing gradients rescaled with regularization term Ω

$$\Omega = \sum_{1 \leq k \leq t} \Omega_k = \sum_{1 \leq k \leq t} \left(\frac{\left\| \frac{\partial \mathcal{E}_t}{\partial \mathbf{h}_k} \frac{\partial \mathbf{h}_k}{\partial \mathbf{h}_{k-1}} \right\|}{\left\| \frac{\partial \mathcal{E}_t}{\partial \mathbf{h}_k} \right\|} - 1 \right)^2$$

- Can result in a 'tug-of-war' scenario where clipping and regularization work diametrically opposed

Experiments and Results

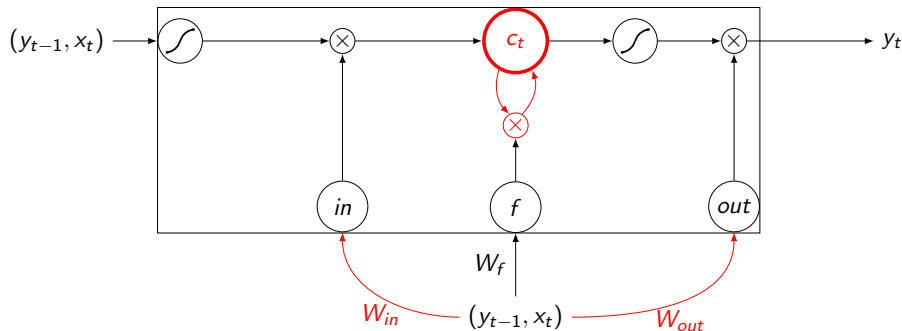
- Sequence with $\{A, B\}$ at beginning and middle
e.g. $(A, \#, \#, \dots, \#, \#, B, \#, \#, \dots, \#, \#) = (A, B)$
- Can generalize to sequences twice as long as the training data



Long Short-Term Memory

- Degenerating gradients are caused by passing the error through many time steps: Can we store the information instead?
- LSTMs store the gradient in a memory cells
- Used nowadays in almost all commercial speech recognition systems
- Come in many different flavors with 'Gated Recurrent Unit' being a very popular one (Reset & Update gate)

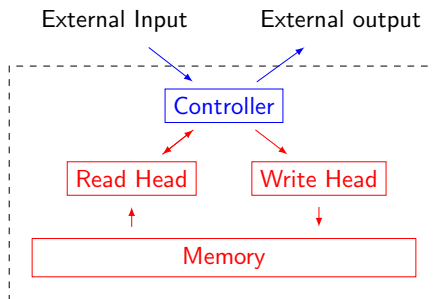
Long Short-Term Memory



He *turned* the radio which was built by Motorola in 1965 *off*.

Neural Turing Machine

- NTM learns to read, write and process information
- Differentiable attention mechanism to interact with memory
- Content and location focusing



Thank You For Your Attention

Sources

- Pascanu, R., Mikolov, T., Bengio, Y. (2013). On the difficulty of training recurrent neural networks. In *Proceedings of the 30th International Conference on Machine Learning*, 1310–1318
- Graves, A. (2013). Generating Sequences With Recurrent Neural Networks, arXiv:1308.0850v5
- Graves, A., Wayne, G., Danihelka, I. (2014). Neural Turing Machines, arXiv:1410.5401v2
- Kumar, A. et al (2016). Ask Me Anything: Dynamic Memory Networks for Natural Language Processing. In *Proceedings of the 33rd International Conference on Machine Learning*, 1378–1387