

Neural Turing Machine

Author: Alex Graves, Greg Wayne, Ivo Danihelka Presented By: Tinghui Wang (Steve)



Introduction

- Neural Turning Machine: Couple a Neural Network with external memory resources
- The combined system is analogous to TM, but differentiable from end to end
- NTM Can infer simple algorithms!

In Computation Theory...

Finite State Machine $(\Sigma, S, s_0, \delta, F)$

- •Σ: Alphabet
- S: Finite set of states
- *s*₀: Initial state
- δ : state transition function: $S \times \Sigma \rightarrow P(S)$
- F: set of Final States



In Computation Theory...

Push Down Automata $(Q, \Sigma, \Gamma, \delta, q_0, Z, F)$

- •Σ: Input Alphabet
- Γ: Stack Alphabet
- Q: Finite set of states
- • q_0 : Initial state
- δ : state transition function: $Q \times \Sigma \times \Gamma \rightarrow P(Q) \times \Gamma$
- F: set of Final States



Turing Machine

 $(\boldsymbol{Q}, \boldsymbol{\Gamma}, \boldsymbol{b}, \boldsymbol{\Sigma}, \boldsymbol{\delta}, \boldsymbol{q}_0, \boldsymbol{F})$

- Q: Finite set of States
- q₀: Initial State
- Γ: Tape alphabet (Finite)
- *b*: Blank Symbol (occurs infinitely on the tape)
- Σ : Input Alphabet ($\Sigma = \Gamma \setminus \{b\}$)
- δ : Transition Function $(Q \setminus F) \times \Gamma \rightarrow Q \times \Gamma \times \{L, R\}$
- F: Final State (Accepted)





Neural Turing Machine - Add Learning to TM!!

Finite State Machine (Program)



Neural Network (I Can Learn!!)



A Little History on Neural Network

- 1950s: Frank Rosenblatt, Perceptron classification based on linear predictor
- 1969: Minsky Proof that Perceptron Sucks! Not able to learn XOR
- 1980s: Back propagation
- 1990s: Recurrent Neural Network, Long Short-Term Memory
- Late 2000s now: Deep Learning, Fast Computer

Feed-Forward Neural Net and Back Propagation



Rumelhart, David E.; Hinton, Geoffrey E.; Williams, Ronald J. (8 October 1986). "Learning representations by back-propagating errors". *Nature* **323** (6088): 533–536.

• Total Input of unit *j*:

$$x_j = \sum_i y_i w_{ji}$$

- Output of unit *j* (logistic function) $y_j = \frac{1}{1 + e^{-x_j}}$
- Total Error (mean square) $E_{total} = \frac{1}{2} \sum_{c} \sum_{j} (d_{c,j} - y_{c,j})^2$
- Gradient Descent $\frac{\partial E}{\partial w_{ji}} = \frac{\partial E}{\partial y_j} \frac{\partial y_j}{\partial x_j} \frac{\partial x_j}{\partial w_{ji}}$

4/6/2015

Tinghui Wang, EECS, WSU

Recurrent Neural Network



R. J. Williams and D. Zipser. Gradient-based learning algorithms for recurrent networks and their computational complexity. In Back-propagation: Theory, Architectures and Applications. Hillsdale, NJ: Erlbaum, 1994.

- y(t): n-tuple of outputs at time t
- $x^{net}(t)$: m-tuple of external inputs at time t
- U: set of indices k that x_k is output of unit in network
- I: set of indices k that x_k is external input

•
$$x_k = \begin{cases} x_k^{net}(t), & if \ k \in I \\ y_k(t), & if \ k \in R \end{cases}$$

• $y_k(t+1) = f_k(s_k(t+1))$ where f_k is a differential squashing function

•
$$s_k(t+1) = \sum_{i \in I \cup U} w_{ki} x_i(t)$$

Recurrent Neural Network – Time Unrolling



• Error Function:

$$J(t) = -\frac{1}{2} \sum_{k \in U} |e_k(t)|^2 = -\frac{1}{2} \sum_{k \in U} |d_k(t) - y_k(t)|^2$$
$$J^{Total}(t', t) = \sum_{\tau=t'+1}^t J(\tau)$$

• Back Propagation through Time: $\varepsilon_k(t) = \frac{\partial J(t)}{\partial y_k(t)} = e_k(t) = d_k(t) - y_k(t)$

$$\delta_{k}(\tau) = \frac{\partial J(t)}{\partial s_{k}(\tau)} = \frac{\partial J(t)}{\partial y_{k}(\tau)} \frac{\partial y_{k}(\tau)}{\partial s_{k}(\tau)} = f_{k}'(s_{k}(\tau))\varepsilon_{k}(\tau)$$

$$\varepsilon_{k}(\tau-1) = \frac{\partial J(t)}{\partial y_{k}(\tau-1)} = \sum_{l \in U} \frac{\partial J(t)}{\partial s_{l}(\tau)} \frac{\partial s_{l}(\tau)}{\partial y_{k}(\tau-1)} = \sum_{l \in U} w_{lk}\delta_{l}(\tau)$$

$$\frac{\partial J^{Total}}{\partial w_{ji}} = \sum_{\tau=t'+1}^{t} \frac{\partial J(t)}{\partial s_{j}(\tau)} \frac{\partial s_{j}(\tau)}{\partial w_{ji}}$$

Tinghui Wang, EECS, WSU

RNN & Long Short-Term Memory

- RNN is Turing Complete
 - With Proper Configuration, it can simulate any sequence generated by a Turing Machine
- Error Vanishing and Exploding
 - Consider single unit self loop:

$$\delta_{k}(\tau) = \frac{\partial J(t)}{\partial s_{k}(\tau)} = \frac{\partial J(t)}{\partial y_{k}(\tau)} \frac{\partial y_{k}(\tau)}{\partial s_{k}(\tau)} = f_{k}'(s_{k}(\tau))\varepsilon_{k}(\tau) = f_{k}'(s_{k}(\tau))w_{kk}\delta_{k}(\tau+1)$$

Long Short-term Memory

 $\begin{array}{c|c} net_{c_{j}} & s_{c_{j}} = s_{c_{j}} + g y^{in_{j}} \\ \hline g & g y^{in_{j}} & 1.0 \\ \hline W_{c_{j}i} & y^{in_{j}} \\ \hline W_{in_{j}i} & y^{out_{j}} \\ \hline W_{in_{j}i} & y^{out_{j}} \\ \hline W_{out_{j}i} & y^{out_{j}} \\ \hline \end{array} \begin{array}{c} y^{out_{j}} \\ \hline y^{out_{j}} \\ \hline y^{out_{j}} \\ \hline y^{out_{j}} \\ \hline \end{array} \begin{array}{c} y^{out_{j}} \\ \hline y^{out_{j}} \\ \hline \end{array} \begin{array}{c} net_{out_{j}} \\ \hline y^{out_{j}} \\ \hline \end{array} \end{array}$

Hochreiter, S. and Schmidhuber, J. (1997). Long shortterm memory. Neural computation, 9(8):1735–1780.

Purpose of Neural Turing Machine

 Enrich the capabilities of standard RNN to simplify the solution of algorithmic tasks



NTM Read/Write Operation



 $w_t(i)$: vector of weights over the N locations emitted by a read head at time t e_t : Erase Vector, emitted by the header a_t : Add Vector, emitted by the header w_t : weights vector

Attention & Focus

 Attention and Focus is adjusted by weights vector across whole memory bank!



Content Based Addressing

• Find Similar Data in memory



May not depend on content only...

Gating Content Addressing



Convolutional Shift – Location Addressing



Convolutional Shift – Location Addressing

 γ_t : sharpening factor



Neural Turing Machine - Experiments

- Goal: to demonstrate NTM is
 - Able to solve the problems
 - By learning compact internal programs
- Three architectures:
 - NTM with a feedforward controller
 - NTM with an LSTM controller
 - Standard LSTM controller
- Applications
 - Copy
 - Repeat Copy
 - Associative Recall
 - Dynamic N-Grams
 - Priority Sort

NTM Experiments: Copy

- Task: Store and recall a long sequence of arbitrary information
- Training: 8-bit random vectors with length 1-20
- No inputs while it generates the targets





Tinghui Wang, EECS, WSU

NTM Experiments: Copy – Learning Curve



NTM Experiments: Repeated Copy

- Task: Repeat a sequence a specified number of times
- Training: random length sequences of random binary vectors followed by a scalar value indicating the desired number of copies





Tinghui Wang, EECS, WSU

NTM Experiments: Repeated Copy – Learning Curve



Tinghui Wang, EECS, WSU

NTM Experiments: Associative Recall

- Task: ask the network to produce the next item, given current item after propagating a sequence to network
- Training: each item is composed of three six-bit binary vectors, 2-8 items every episode





NTM Experiments: Dynamic N-Grams

- Task: Learn N-Gram Model rapidly adapt to new predictive distribution
- Training: 6-Gram distributions over binary sequences. 200 successive bits using look-up table by drawing 32 probabilities from Beta(.5,.5) distribution
- Compare to Optimal Estimator: $P(B = 1 | N_1, N_2, c) = \frac{N_1 + 0.5}{N_1 + N_0 + 1}$





NTM Experiments: Priority Sorting

- Task: Sort Binary Vector based on priority
- Training: 20 binary vectors with corresponding priorities, output 16 highest-priority vectors





Tinghui Wang, EECS, WSU

NTM Experiments: Learning Curve



Figure 13: Dynamic N-Gram Learning Curves.

Figure 18: Priority Sort Learning Curves.

Some Detailed Parameters

NTM with						
Feed Forward						
Neural Network						

NTM with LSTM

LOIN

Task	#Heads	Controller Size	Memory Size	Learning Rate	#Parameters
Сору	1	100	128×20	10^{-4}	17,162
Repeat Copy	1	100	128 imes 20	10^{-4}	16,712
Associative	4	256	128×20	10^{-4}	146,845
N-Grams	1	100	128 imes 20	$3 imes 10^{-5}$	14,656
Priority Sort	8	512	128×20	$3 imes 10^{-5}$	508,305
Task	#Heads	Controller Size	Memory Size	Learning Rate	#Parameters
Сору	1	100	128 imes 20	10^{-4}	67, 561
Repeat Copy	1	100	128×20	10^{-4}	66, 111
Associative	1	100	128×20	10^{-4}	70,330
N-Grams	1	100	128×20	$3 imes 10^{-5}$	61,749
Priority Sort	5	2×100	128×20	$3 imes 10^{-5}$	269,038
Task	Network Size Learning Rate #Parameters				
Сору	3 imes 2	$256 3 \times 10$	$^{-5}$ 1,352,	969	

 $3 imes 10^{-5}$

 10^{-4}

 10^{-4}

 3×10^{-5}

5,312,007

1,344,518

331,905

384,424

 3×512

 3×256

 3×128

 3×128

Repeat Copy

Associative

Priority Sort

N-Grams

Additional Information

- Literature Reference:
 - Rumelhart, David E.; Hinton, Geoffrey E.; Williams, Ronald J. (8 October 1986). "Learning representations by back-propagating errors". *Nature* 323 (6088): 533–536
 - R. J. Williams and D. Zipser. Gradient-based learning algorithms for recurrent networks and their computational complexity. In Back-propagation: Theory, Architectures and Applications. Hillsdale, NJ: Erlbaum, 1994.
 - Hochreiter, S. and Schmidhuber, J. (1997). Long short-term memory. Neural computation, 9(8):1735–1780
 - Y. Bengio and Y. LeCun, Scaling Learning Algorithms towards AI, Large-Scale Kernel Machines, Bottou, ., Chapelle, O., DeCoste, D., and Weston, J., Eds., MIT Press, 2007.
- NTM Implementation
 - https://github.com/shawntan/neural-turing-machines
- NTM Reddit Discussion
 - <u>https://www.reddit.com/r/MachineLearning/comments/2m9mga/implementation of neural turing machines/</u>

