



# Implicitly-Defined Neural Networks for Sequence Labeling

WeiYang

[weiyang@godweiyang.com](mailto:weiyang@godweiyang.com)

[www.godweiyang.com](http://www.godweiyang.com)

East China Normal University  
Department of Computer Science

2017.11.03



# Outline

Outline

Traditional RNN

INN

Experiments





# Formula

- input sequence

$$[\xi_1, \xi_2, \dots, \xi_n] \quad (1)$$

- states production

$$h_1 = f(\xi_1, h_s)$$

$$h_2 = f(\xi_2, h_1)$$

...

$$h_n = f(\xi_n, h_{n-1})$$

(2)

- LSTM & GRU



# Structure

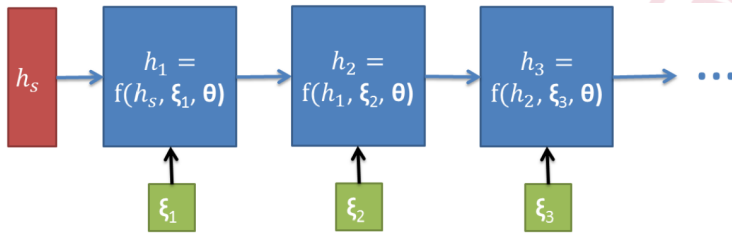


Figure 1: Traditional RNN structure.

# Formula

- implicit hidden layer

where

$$H = [h_1, h_2, \dots, h_n] \quad (3)$$

$$h_1 = f(h_s, h_2, \xi_1)$$

...

$$h_i = f(h_{i-1}, h_{i+1}, \xi_i) \quad (4)$$

...

$$h_n = f(h_{n-1}, h_e, \xi_n)$$

- INN

$$\xi = g(\theta, X)$$

$$H = F(\theta, \xi, H) \quad (5)$$

$$L = \ell(\theta, H, Y)$$



# Structure

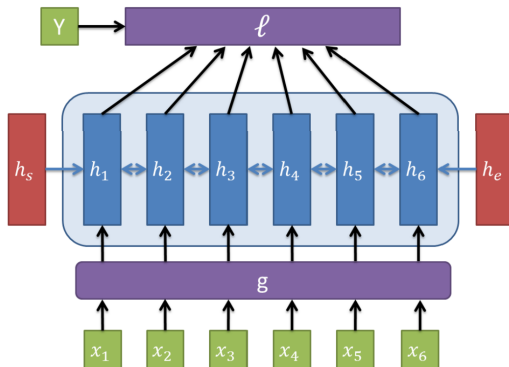


Figure 2: Proposed INN Architecture



# Forward propagation

- solve the equation

$$H = F(H) \quad (6)$$

- approximate Newton solve

$$H_{n+1} = H_n - (I - \nabla_H F)^{-1}(H_n - F(H_n)) \quad (7)$$

- Krylov subspace methods(BiCG-STAB method)



# Backward propagation

- gradient of the loss function

$$\nabla_{\theta} L = \nabla_{\theta} \ell + \nabla_H \ell \nabla_{\theta} H \quad (8)$$

where

$$\nabla_{\theta} H = (I - \nabla_H F)^{-1} (\nabla_{\theta} F + \nabla_{\xi} F \nabla_{\theta} \xi) \quad (9)$$

so

$$\nabla_{\theta} L = \nabla_{\theta} \ell + \nabla_H \ell (I - \nabla_H F)^{-1} (\nabla_{\theta} F + \nabla_{\xi} F \nabla_{\theta} \xi) \quad (10)$$

- Krylov subspace methods (BiCG-STAB method)



# Transition Functions

$$\begin{aligned}
 h_t &= (1 - z_t)\hat{h}_t + z_t\tilde{h}_t \\
 \tilde{h}_t &= \tanh(Wx_t + U(r_t\hat{h}_t) + \tilde{b}) \\
 z_t &= \sigma(W_zx_t + U_z\hat{h}_t + b_z) \\
 r_t &= \sigma(W_rx_t + U_r\hat{h}_t + b_r)
 \end{aligned}
 \tag{11}$$

where

$$\begin{aligned}
 \hat{h}_t &= sh_{t-1} + (1 - s)h_{t+1} \\
 s &= \frac{s_p}{s_p + s_n} \\
 s_p &= \sigma(W_px_t + U_ph_{t-1} + b_p) \\
 s_n &= \sigma(W_nx_t + U_nh_{t+1} + b_n)
 \end{aligned}
 \tag{12}$$



## POS

Architecture	WSJ Accuracy
GRU	96.43
LSTM	96.47
Bidirectional GRU	97.28
b-LSTM	97.25
INN	<b>97.37</b>
Stanford POS Tagger	97.33

Table 2: Tagging performance relative to recurrent architectures and Stanford POS Tagger.