

CS 4248

Natural Language Processing

Professor NG Hwee Tou
Department of Computer Science
School of Computing
National University of Singapore
ngh@comp.nus.edu.sg

Materials

- NNM4NLP Chapter 14, 15, 16

Characteristics of Natural Language

- Sequence data
 - Word (a sequence of characters)
 - Sentence (a sequence of words)
 - Document (a sequence of sentences)
- Convolution neural networks
 - Modeling of word order restricted to mostly local patterns (ngrams)

Recurrent Neural Networks (RNN)

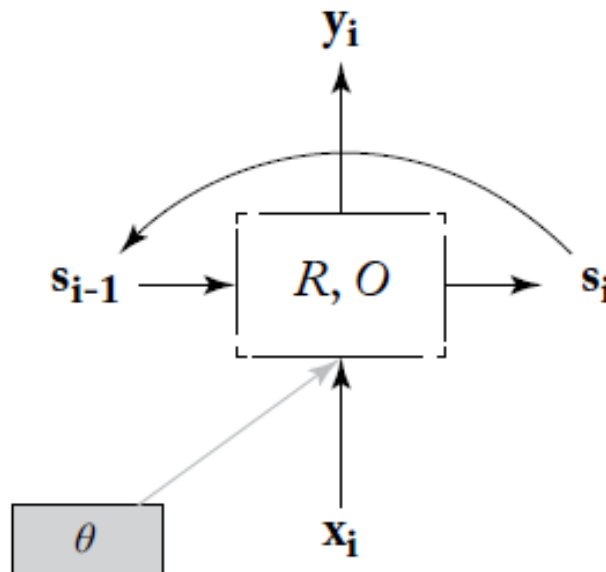
- Jeffrey L. Elman, Finding structure in time. *Cognitive Science*, March 1990
- Capture subtle patterns and regularities in sequences
- Model non-Markovian dependencies
- Consider infinite window and zoom in on informative sequential patterns in the window

Recurrent Neural Networks

- Represent an arbitrarily sized sequence of inputs in a fixed size vector
- Output can be fed into a larger network downstream

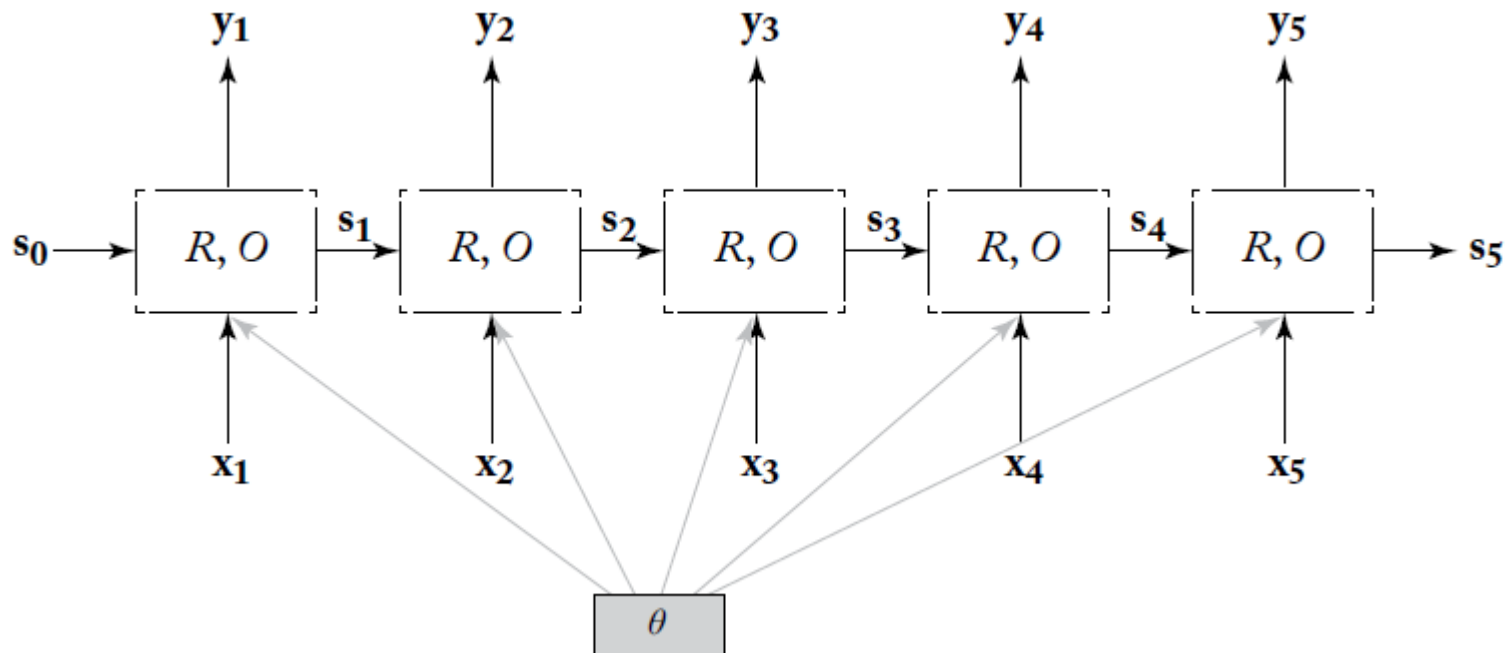
Graphical Representation of an RNN

Recursive view:



Graphical Representation of an RNN

Unrolled view:



- Unrolled view makes clear that RNN is a deep neural network with many layers
- The **same** parameters θ are shared across all time steps

The RNN Abstraction

$$\text{RNN}^*(\mathbf{x}_{1:n}; \mathbf{s}_0) = \mathbf{y}_{1:n}$$

$$\mathbf{s}_i = R(\mathbf{s}_{i-1}, \mathbf{x}_i)$$

$$\mathbf{y}_i = O(\mathbf{s}_i)$$

$$\mathbf{x}_i \in \mathbb{R}^{d_{in}}, \mathbf{y}_i \in \mathbb{R}^{d_{out}}, \mathbf{s}_i \in \mathbb{R}^{f(d_{out})}$$

- The functions R and O are the **same** across the sequence positions
- The last \mathbf{y}_n encodes the entire input sequence and can be used for further prediction

The RNN Abstraction

$$\begin{aligned} s_4 &= R(s_3, x_4) \\ &= R(\overbrace{R(s_2, x_3)}^{s_3}, x_4) \\ &= R(R(\overbrace{R(s_1, x_2)}^{s_2}), x_3), x_4) \\ &= R(R(R(\overbrace{R(s_0, x_1)}^{s_1}), x_2), x_3), x_4) \end{aligned}$$

RNN

- Sensitive to the order of words in a sequence

$$\mathbf{s}_i = R_{\text{RNN}}(\mathbf{s}_{i-1}, \mathbf{x}_i) = g([\mathbf{s}_{i-1}; \mathbf{x}_i]\mathbf{W} + \mathbf{b})$$

$$\mathbf{y}_i = O_{\text{RNN}}(\mathbf{s}_i) = \mathbf{s}_i$$

$$\mathbf{x}_i \in \mathbb{R}^{d_x}, \mathbf{s}_i, \mathbf{y}_i \in \mathbb{R}^{d_s}, \mathbf{b} \in \mathbb{R}^{d_s}$$

$$\mathbf{W} \in \mathbb{R}^{(d_s+d_x) \times d_s}$$

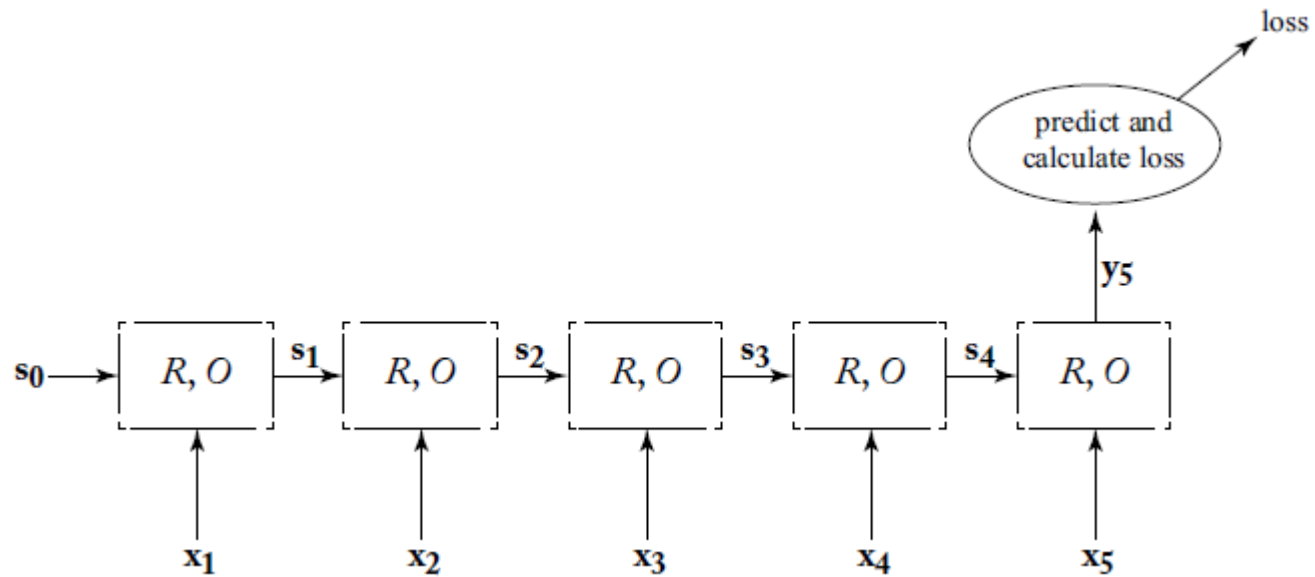
The RNN Abstraction

- E.g., a model for predicting the conditional probability of an event e given the sequence $\mathbf{x}_{1:n}$:
- $p(e = j | \mathbf{x}_{1:n}) = \text{softmax}(\mathbf{y}_n \cdot \mathbf{W} + \mathbf{b})_{[j]}$
- RNN serves as a trainable component in a larger network
- RNN allows conditioning on the entire history $\mathbf{x}_1, \dots, \mathbf{x}_n$ without making the Markov assumption

RNN Training

- Backpropagation through time (BPTT)
 - Create the unrolled computation graph for a given input sequence
 - Add a loss node
 - Use backpropagation to compute the gradients w.r.t. the loss function

RNN as Acceptor



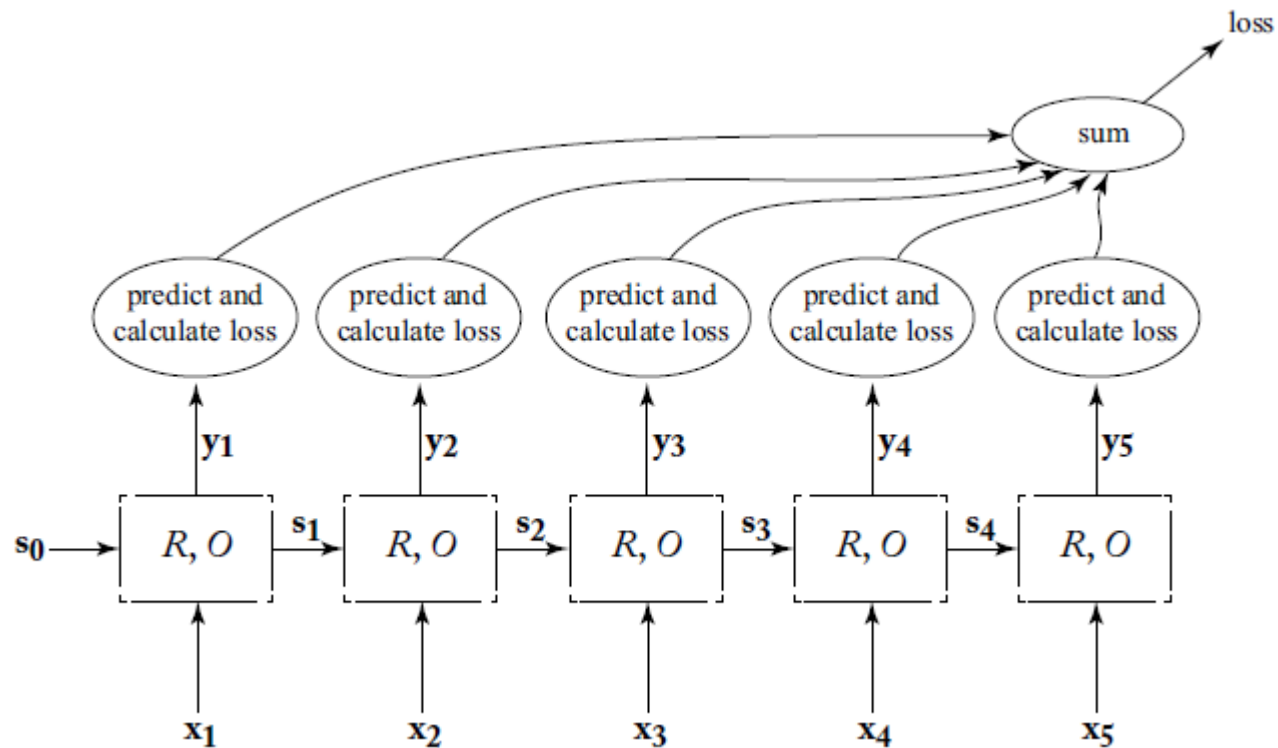
RNN as Acceptor

- Based on the final state y_n , make a prediction
- E.g., read the words of a sentence s and predict whether s conveys a positive or negative sentiment

RNN as Encoder

- y_n is treated as an encoding of x_1, \dots, x_n , and used together with other features
- E.g., an extractive document summarizer first uses an RNN to encode an input document with y_n , then uses y_n together with other features to select the sentences to include in the summary

RNN as Transducer

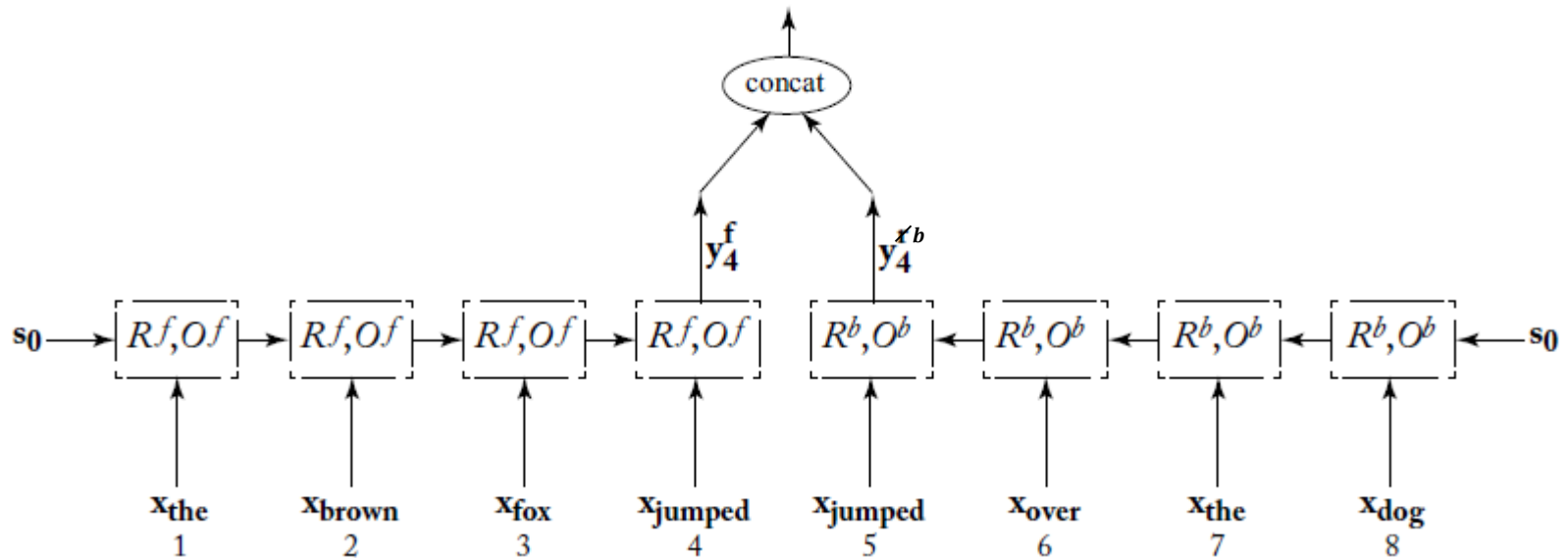


RNN as Transducer

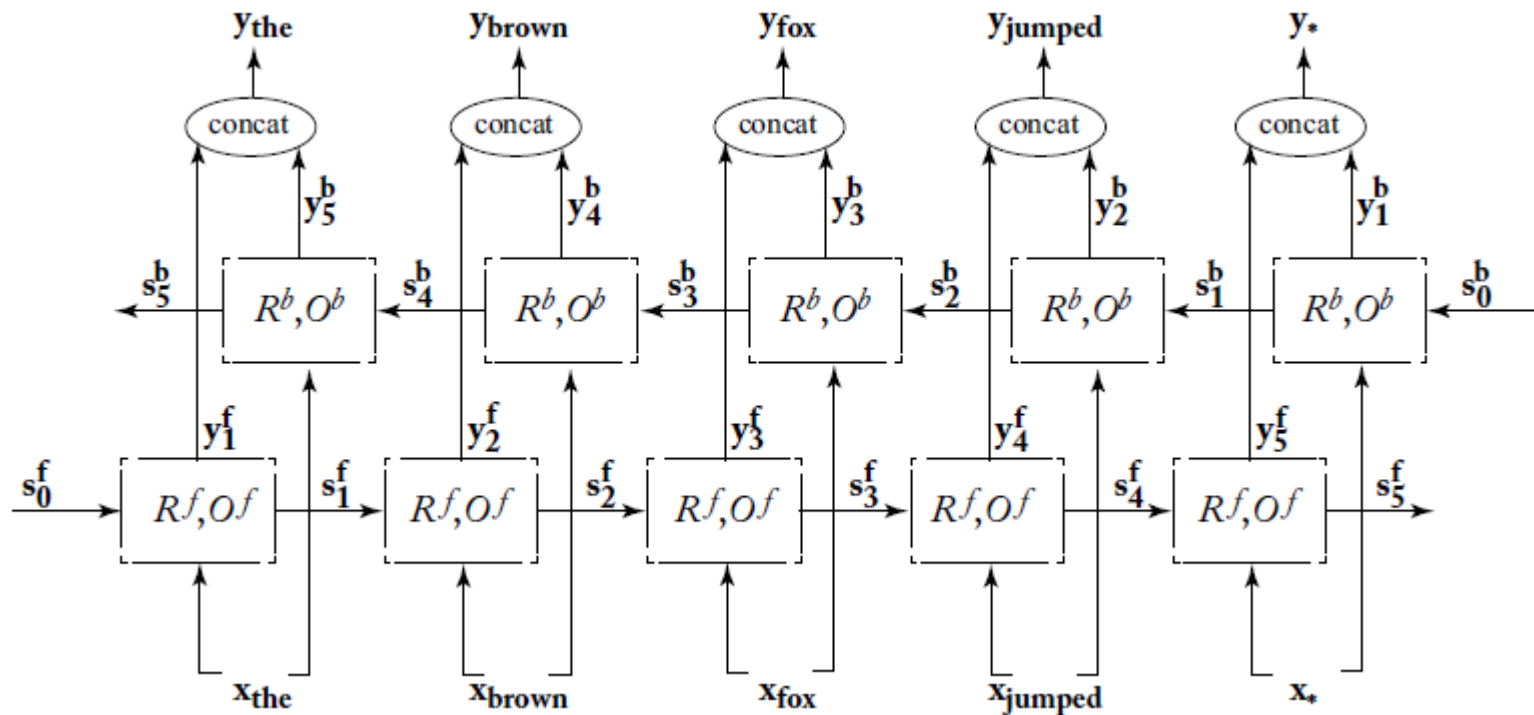
- Produce an output \hat{t}_i after each input x_i is read
- Local loss $L_{\text{local}}(\hat{t}_i, t_i)$
- $L(\hat{t}_{1:n}, t_{1:n}) = \sum_{i=1}^n L_{\text{local}}(\hat{t}_i, t_i)$
- E.g., a sequence tagger that predicts the tag of each word
- E.g., language modeling predicts the next word based on x_1, \dots, x_i

Bidirectional RNN

- Bidirectional RNN (biRNN): use **both** the previous and following words (context)



Bidirectional RNN



Bidirectional RNN

- Input sequence: $x_{1:n}$
- 2 separate states for each position i : s_i^f and s_i^b
- s_i^f : based on x_1, x_2, \dots, x_i
- s_i^b : based on x_n, x_{n-1}, \dots, x_i

Output vector at position i :

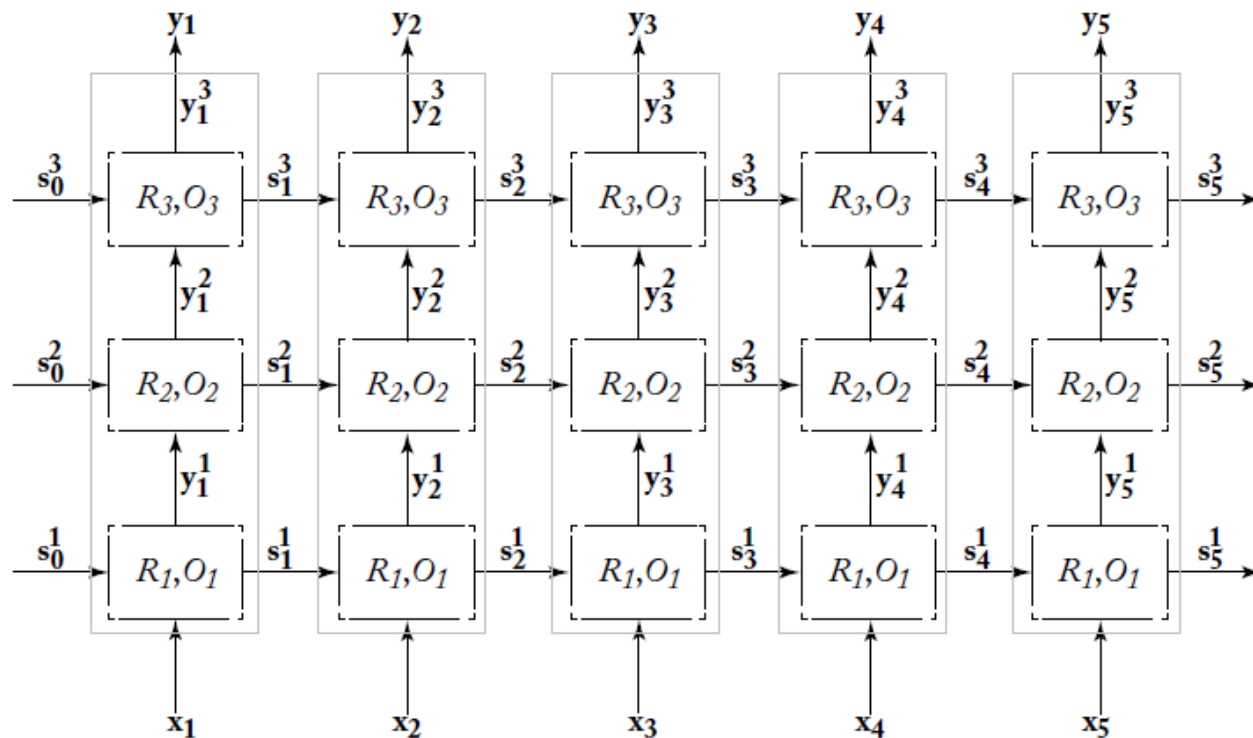
- $\text{biRNN}(x_{1:n}, i) = y_i = [\text{RNN}^f(x_{1:i}); \text{RNN}^b(x_{n:i})] = [y_i^f; y_i^b] = [O^f(s_i^f); O^b(s_i^b)]$

Bidirectional RNN

- $\text{biRNN}^*(\mathbf{x}_{1:n}) = \mathbf{y}_{1:n} =$
 $\text{biRNN}(\mathbf{x}_{1:n}, 1), \dots, \text{biRNN}(\mathbf{x}_{1:n}, n)$
- biRNN is very effective for sequence tagging tasks
- biRNN is used as a trainable feature extractor

Multi-Layer (Stacked) RNN

Deep RNN



Multi-Layer (Stacked) RNN

- k RNNs, the j th RNN has states $s_{1:n}^j$ and outputs $y_{1:n}^j$
- Input for the first RNN: $x_{1:n}$
- Input for the j th RNN ($j \geq 2$) = Output of the RNN below it, $y_{1:n}^{j-1}$
- Output of the entire deep RNN = Output of the last RNN $y_{1:n}^k$
- Deep RNNs outperform shallower ones on some tasks (e.g., neural machine translation)

CBOW

- CBOW can be considered as an RNN without taking into account order of words

$$\begin{aligned} \mathbf{s}_i &= R_{\text{CBOW}}(\mathbf{s}_{i-1}, \mathbf{x}_i) = \mathbf{s}_{i-1} + \mathbf{x}_i \\ \mathbf{y}_i &= O_{\text{CBOW}}(\mathbf{s}_i) = \mathbf{s}_i \\ \mathbf{x}_i, \mathbf{y}_i, \mathbf{s}_i &\in \mathbb{R}^{d_s} \end{aligned}$$

- It follows that $\mathbf{s}_n = \sum_{i=1}^n \mathbf{x}_i$ (assume $\mathbf{s}_0 = \mathbf{0}$)

Modeling with RNNs

- RNN as feature extractor
- RNN as acceptor

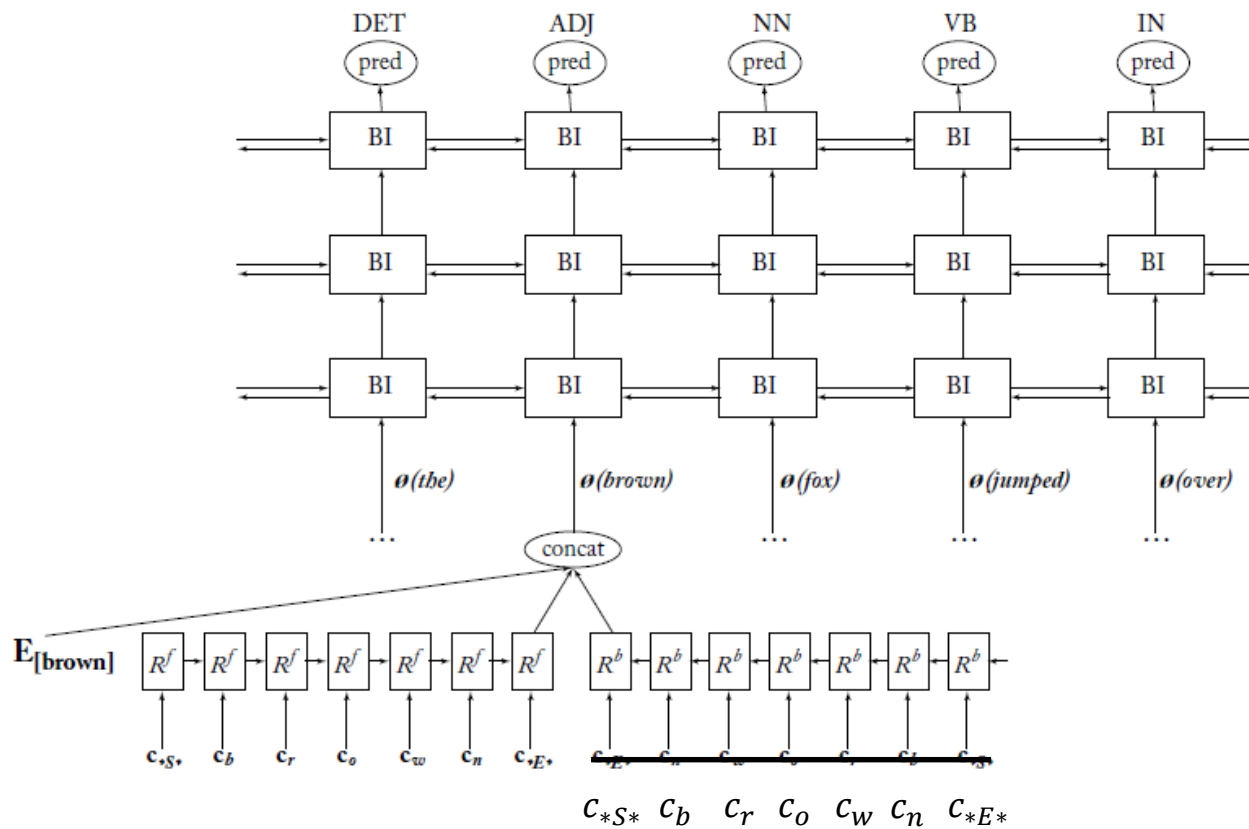
RNN as Feature Extractor

- RNN is used as a trainable feature extractor, to replace parts of the more traditional feature extractor
- E.g., POS tagging

POS Tagging

- A sequence tagging task: assign a tag to each of the n input words
- biRNN

POS Tagging



POS Tagging

- s : input sentence
- $w_i = c_1, \dots, c_l$
- c_i : character embedding vector for c_i
- $\mathbf{x}_i = \phi(s, i) = [\mathbf{E}_{[w_i]}; \text{RNN}^f(\mathbf{c}_{1:l}); \text{RNN}^b(\mathbf{c}_{l:1})]$
- $p(t_i = j | w_1, \dots, w_n) = \text{softmax}(\text{MLP}(\text{biRNN}(\mathbf{x}_{1:n}, i)))_{[j]}$
- Trained with cross-entropy loss function

POS Tagging

- Condition on previous POS tags predicted
- Condition on k previous POS tags predicted using tag embeddings $E_{[t]}$
- $p(t_i = j | w_1, \dots, w_n, t_{i-1}, \dots, t_{i-k}) = \text{softmax}(\text{MLP}([\text{biRNN}(\mathbf{x}_{1:n}, i); E_{[t_{i-1}]}; \dots; E_{[t_{i-k}]}]))_{[j]}$
- Condition on the entire sequence of POS tags predicted
- $p(t_i = j | w_1, \dots, w_n, t_{1:i-1}) = \text{softmax}(\text{MLP}([\text{biRNN}(\mathbf{x}_{1:n}, i); \text{RNN}^t(\mathbf{t}_{1:i-1})]))_{[j]}$

POS Tagging

- Character embeddings take into account suffix, prefix, and orthographic (capitalization, hyphens, digits) cues important for POS tagging
- POS tag prediction is conditioned on the entire input sentence
- biRNN is used as feature extractor, with no manually defined features

RNN as Acceptor

- Read an input sequence and produce a binary or multi-class answer
- E.g., sentiment classification

Sentence-Level Sentiment Classification

- Given a sentence s , classify s as positive or negative
- Examples:
 - Positive: *It's not life-affirming – it's vulgar and mean, but I liked it.*
 - Negative: *It's disappointing that it only manages to be decent instead of dead brilliant.*
- Need to deal with negation, sarcasm, sentence structure, etc.

Sentence-Level Sentiment Classification

- $p(\text{label} = k | w_{1:n}) = \hat{y}_{[k]} \quad k = 1, 2$
- $\hat{y} = \text{softmax}(\text{MLP}(\text{RNN}(\mathbf{x}_{1:n})))$
- $\mathbf{x}_{1:n} = \mathbf{E}_{[w_1]}, \dots, \mathbf{E}_{[w_n]}$
- Trained with cross-entropy loss function
- Word embedding matrix \mathbf{E} is initialized using pre-trained word embeddings (e.g., from Word2Vec)

Sentence-Level Sentiment Classification

- Bidirectional RNN
- $p(\text{label} = k | w_{1:n}) = \hat{y}_{[k]} \quad k = 1, 2$
- $\hat{y} = \text{softmax}(\text{MLP}([\text{RNN}^f(x_{1:n}); \text{RNN}^b(x_{n:1})]))$
- $x_{1:n} = E_{[w_1]}, \dots, E_{[w_n]}$