

DSME 6635: Artificial Intelligence for Business Research

## Deep-Learning-based NLP: RNN and Seq2Seq

Renyu (Philip) Zhang

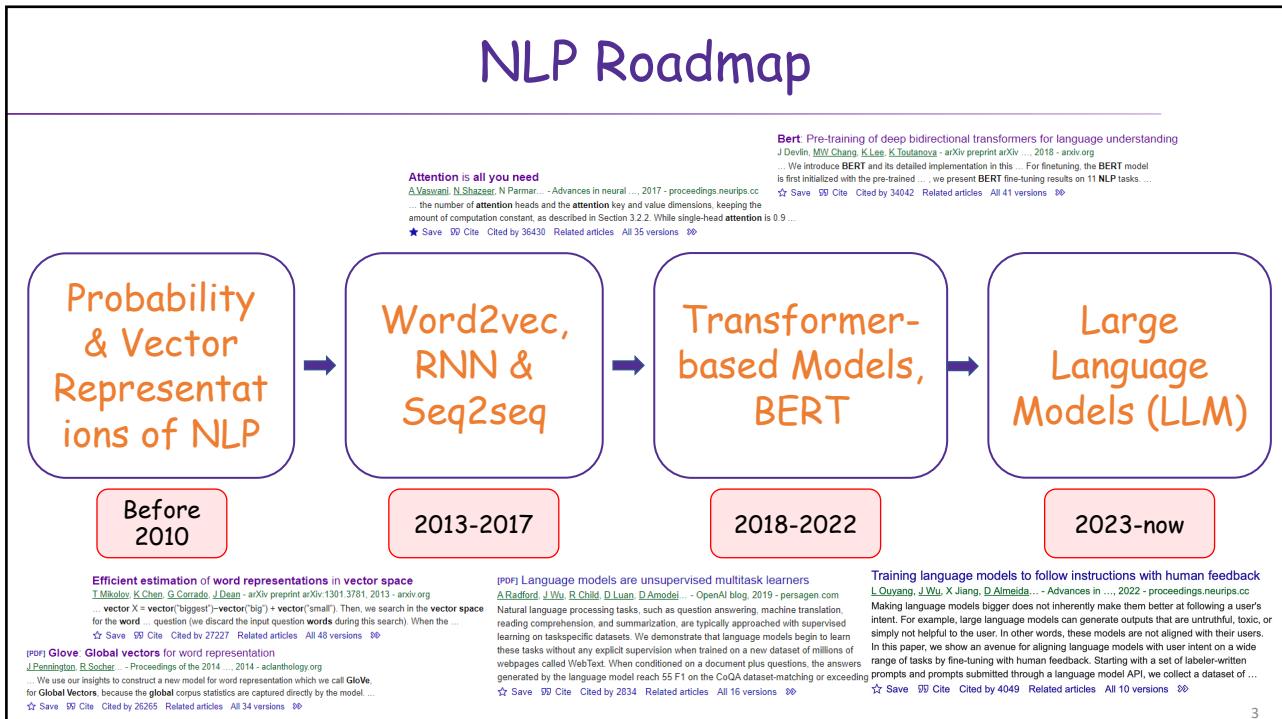
1

## Agenda

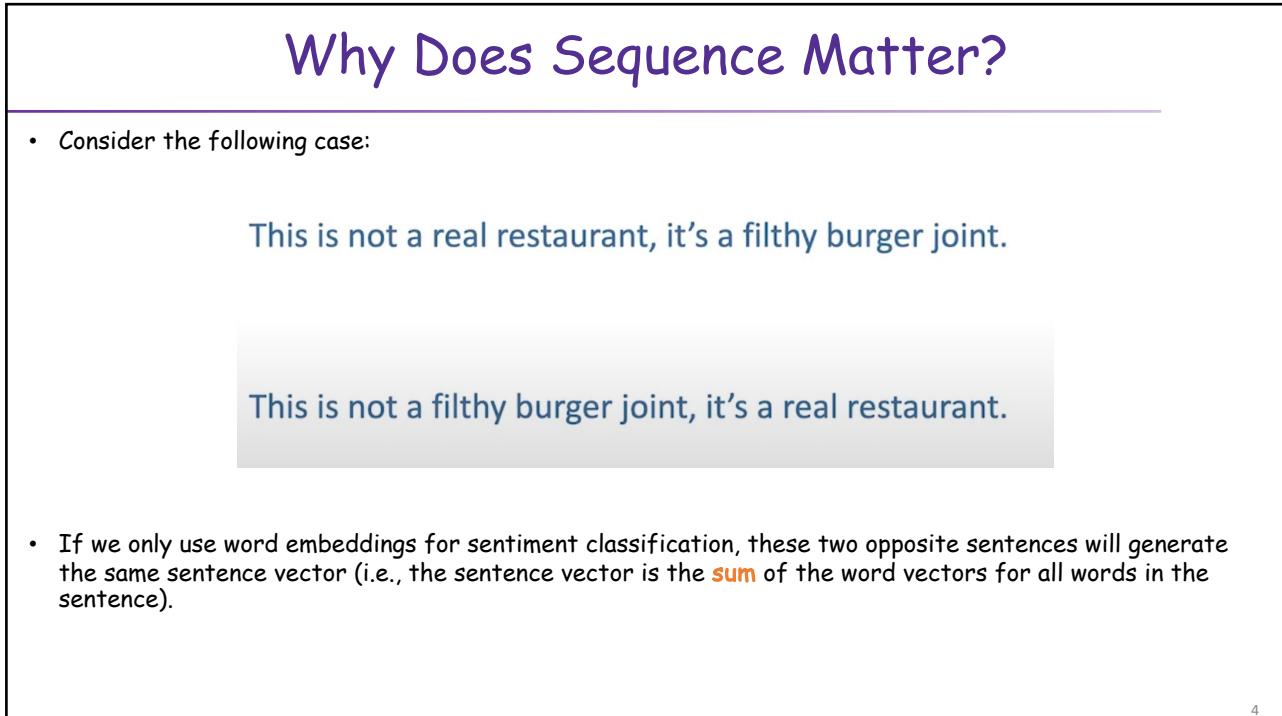
- Vanilla Recurrent Neural Nets (RNN)
- Long Short-Term Memory (LSTM)
- Sequence-to-sequence (Seq2seq)

2

2



3



4

## Back to N-Gram Models

- N-gram model is a **language model** that limits the dependencies on history, a.k.a. **Markov Chain**.

$$P(\mathbf{x}^{(t+1)} | \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)}) = P(\mathbf{x}^{(t+1)} | \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(t-n+2)}) \quad (\text{assumption})$$

*n-1 words*

$$\begin{aligned} \text{prob of a n-gram} &= P(\mathbf{x}^{(t+1)}, \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(t-n+2)}) \\ \text{prob of a (n-1)-gram} &= P(\mathbf{x}^{(t)}, \dots, \mathbf{x}^{(t-n+2)}) \end{aligned} \quad (\text{definition of conditional prob})$$

- Question:** How do we get these  $n$ -gram and  $(n-1)$ -gram probabilities?
- Answer:** By **counting** them in some large corpus of text!

$$\approx \frac{\text{count}(\mathbf{x}^{(t+1)}, \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(t-n+2)})}{\text{count}(\mathbf{x}^{(t)}, \dots, \mathbf{x}^{(t-n+2)})} \quad (\text{statistical approximation})$$

Two important issues:

- Sparsity** (partially addressed by **smoothing**)
- Model Size** ( $n$  is no more than 5, usually **2 or 3**)

- N-gram models that "work" in the era of LLM: <https://arxiv.org/abs/2401.17377>

5

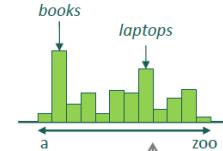
## Recurrent Neural Network (RNN)

Reference: Stanford CS224N, Lecture 5:  
<https://web.stanford.edu/class/cs224n/slides/cs224n-2024-lecture05-rnnlm.pdf>

output distribution

$$\hat{y}^{(t)} = \text{softmax}(\mathbf{U} \mathbf{h}^{(t)} + \mathbf{b}_2) \in \mathbb{R}^{|V|}$$

$$\hat{y}^{(4)} = P(\mathbf{x}^{(5)} | \text{the students opened their})$$



hidden states

$$\mathbf{h}^{(t)} = \sigma(\mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_e \mathbf{e}^{(t)} + \mathbf{b}_1)$$

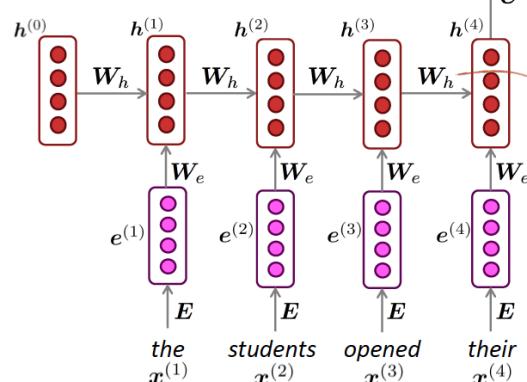
$\mathbf{h}^{(0)}$  is the initial hidden state

The same  $\mathbf{W}_h$  is applied throughout, so it can handle any input sequence length.

word embeddings

$$\mathbf{e}^{(t)} = \mathbf{E} \mathbf{x}^{(t)}$$

words / one-hot vectors  
 $\mathbf{x}^{(t)} \in \mathbb{R}^{|V|}$



6

6

# Recurrent Neural Network (RNN)

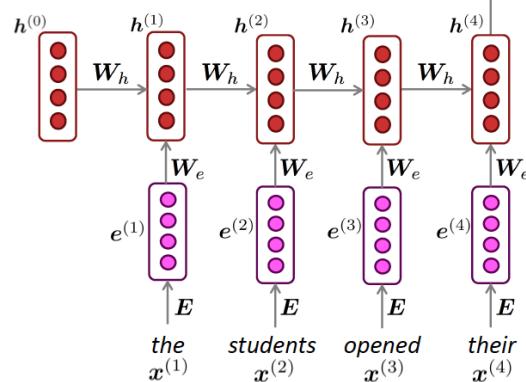
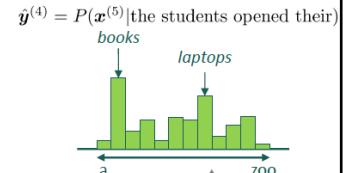
Reference: Stanford CS224N, Lecture 5:  
<https://web.stanford.edu/class/cs224n/slides/cs224n-2024-lecture05-rnnlm.pdf>

## RNN Advantages:

- Can process **any input sequence length**.
- Computations (in theory) use **information from many steps back**.
- Same weights are applied to every step, so there's **time-symmetry/invariance** in how inputs are processed.

## RNN Disadvantages:

- Recurrent computations are **slow**.
- In practice, it is challenging to access information from **many steps back**.



7

7

# Training RNN

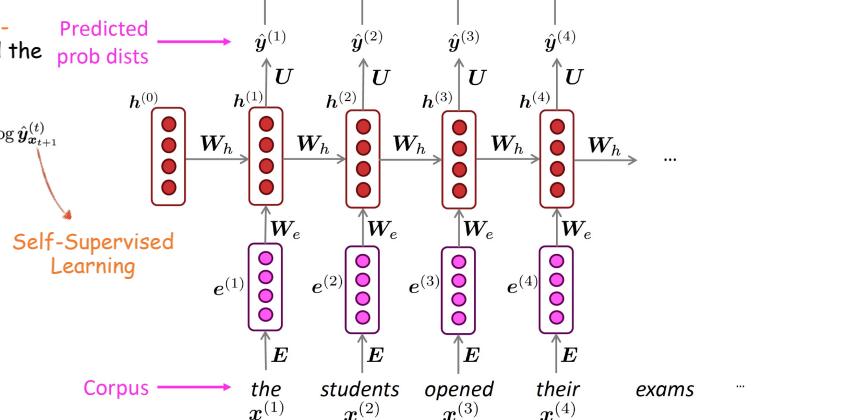
Reference: Stanford CS224N, Lecture 5:  
<https://web.stanford.edu/class/cs224n/slides/cs224n-2024-lecture05-rnnlm.pdf>

- Loss functions in step  $t$  is the **cross-entropy** between the true 1-hot and the predicted prob dists:

$$J^{(t)}(\theta) = CE(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) = - \sum_{w \in V} \mathbf{y}_w^{(t)} \log \hat{y}_w^{(t)} = - \log \hat{y}_{x_{t+1}}^{(t)}$$

- So, the overall loss for the entire training corpus is:

$$J(\theta) = \frac{1}{T} \sum_{t=1}^T J^{(t)}(\theta) = \frac{1}{T} \sum_{t=1}^T -\log \hat{y}_{x_{t+1}}^{(t)}$$



8

8

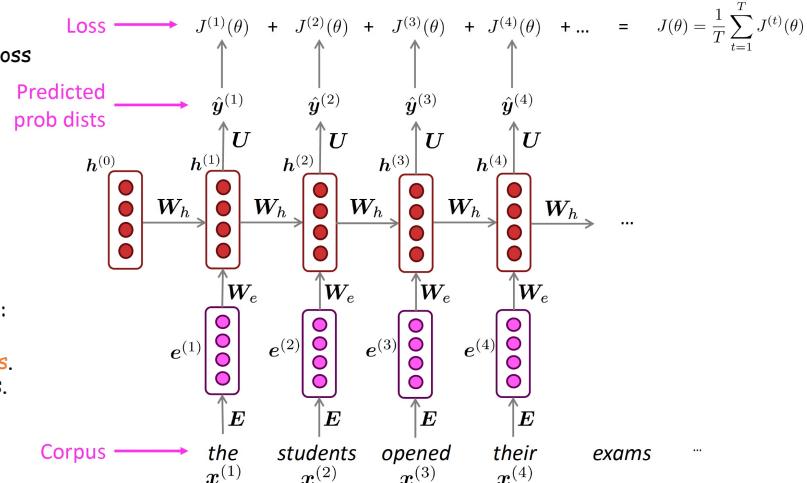
## Training RNN

Reference: Stanford CS224N, Lecture 5:

<https://web.stanford.edu/class/cs224n/slides/cs224n-2024-lecture05-rnnlm.pdf>

- Computing the loss and the gradients across the entire corpus is computationally too expensive.

- In practice, we leverage the idea of SGD: Compute loss and gradients, and update weights with batches of words/sentences. Then repeat on a new batch of sentences.



9

9

## Vanishing (and Exploding) Gradient in RNN

Reference: Stanford CS224N, Lecture 5:

<https://web.stanford.edu/class/cs224n/slides/cs224n-2024-lecture05-rnnlm.pdf>

Backpropagation through time

$$\frac{\partial L}{\partial W_h} \propto \sum_{1 \leq k \leq t} \left( \prod_{l \geq i > k} \frac{\partial \hat{y}_i}{\partial h_{i-1}} \right) \frac{\partial \hat{y}_k}{\partial W_h}$$

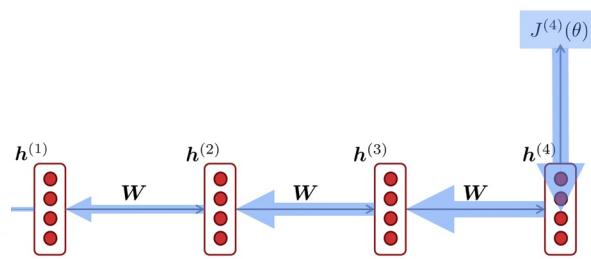
Contribution of hidden state k

Length of the product proportional to how far k is from t

$$\frac{\partial \hat{y}_1}{\partial W_h} \frac{\partial \hat{y}_{t-1}}{\partial W_h} \frac{\partial \hat{y}_{t-2}}{\partial W_h} \frac{\partial \hat{y}_{t-3}}{\partial W_h} \frac{\partial \hat{y}_{t-4}}{\partial W_h} \frac{\partial \hat{y}_{t-5}}{\partial W_h} \frac{\partial \hat{y}_{t-6}}{\partial W_h} \frac{\partial \hat{y}_{t-7}}{\partial W_h} \frac{\partial \hat{y}_{t-8}}{\partial W_h} \frac{\partial \hat{y}_{t-9}}{\partial W_h} \frac{\partial \hat{y}_{t-10}}{\partial W_h}$$

Contribution of hidden state t-10

Vanishing vs. Exploding Gradient



$$\frac{\partial J^{(4)}}{\partial h^{(1)}} = \frac{\partial h^{(2)}}{\partial h^{(1)}} \times \frac{\partial h^{(3)}}{\partial h^{(2)}} \times \frac{\partial h^{(4)}}{\partial h^{(3)}} \times \frac{\partial J^{(4)}}{\partial h^{(4)}}$$

What happens if these are small?

**Vanishing gradient problem:**  
When these are small, the gradient signal gets smaller and smaller as it backpropagates further

Gradient signals from far away will be lost!  
The weights  $W_h$  only capture near effects.

On the difficulty of training recurrent neural networks

R Pascanu, T Mikolov, Y Bengio

International conference on machine learning, 2013 · proceedings.mlr.press

Abstract

There are two widely known issues with properly training recurrent neural networks, the vanishing and the exploding gradient problems detailed in Bengio et al.(1994). In this paper we attempt to improve the understanding of the underlying issues by exploring these problems from an analytical, a geometric and a dynamical systems perspective. Our analysis is used to justify a simple yet effective solution. We propose a gradient norm clipping strategy to deal with exploding gradients and a soft constraint for the vanishing

SHOW MORE ▾

☆ Save 59 Cite Cited by 6901 Related articles All 11 versions ⓘ

10

10

# Gradient Clipping and Skip Connection

Reference: Stanford CS224N, Lecture 5:  
<https://web.stanford.edu/class/cs224n/slides/cs224n-2024-lecture05-rnnlm.pdf>

- The exploding gradient problem is relatively easier to address: **Gradient Clipping**.

- Intuition: Take a **smaller step** in the **same direction**.

- One idea to address vanishing gradient is to create **direct** and **linear pass-through connections** in the model: **Residual/skip connections, attention, etc.**

---

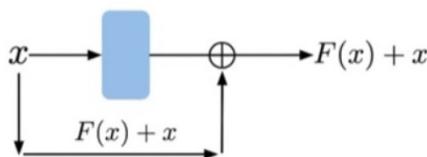
**Algorithm 1** Pseudo-code for norm clipping

```

 $\hat{g} \leftarrow \frac{\partial E}{\partial \theta}$ 
if  $\|\hat{g}\| \geq threshold$  then
     $\hat{g} \leftarrow \frac{threshold}{\|\hat{g}\|} \hat{g}$ 
end if

```

---



Deep residual learning for image recognition

K He, X Zhang, S Ren, J Sun - ... and pattern recognition, 2016 - openaccess.thecvf.com

... Deeper neural networks are more difficult to train. We present a residual learning framework to ease the training of networks that are substantially deeper than those used previously. ...

★ 保存 引用 被引用次数 : 196717 相关文章 所有 76 个版本 ☰

11

11

# Agenda

- Vanilla Recurrent Neural Nets (RNN)
- Long Short-Term Memory (LSTM)
- Sequence-to-sequence (Seq2seq)

12

12

# Long Short-Term Memory RNN (LSTM)

Reference: Stanford CS224N, Lecture 6: <https://web.stanford.edu/class/cs224n/slides/cs224n-2024-lecture06-fancy-rnn.pdf>

- LSTM is a very commonly used RNN architecture that solves the vanishing gradient problem by **disregarding the irrelevant past information** based on the current information.
- LSTM was the **dominant approach** for most NLP tasks in 2013-2015.
- Very interesting history:**
  - Everyone cites Hochreiter and Schmidhuber (1997), but the crucial part of modern LSTM comes from Gers et al. (2000).
  - Recognized as promising only after Graves et al. (2006) which invented CTC (connectionist temporal classification) for speech recognition.
  - Became well-known after Hinton brought LSTM to Google in 2013 (Graves was a student of Schmidhuber and a post-doc of Hinton).

## Long short-term memory

S Hochreiter, J Schmidhuber - Neural computation, 1997 - ieeexplore.ieee.org

... (**short-term memory**, as opposed to **long-term memory**) ... learning what to put in **shortterm memory**, however, take too ... and corresponding teacher signals are **long**. Although theoretically ...

☆ Save 99 Cite Cited by 98734 Related articles All 45 versions ☺

## Learning to forget: Continual prediction with LSTM

F A Gers, J Schmidhuber, F Cummins - Neural computation, 2000 - ieeexplore.ieee.org

Long short-term memory (LSTM; Hochreiter & Schmidhuber, 1997) can solve numerous tasks not solvable by previous **learning** algorithms for recurrent neural networks (RNNs). We ...

☆ Save 99 Cite Cited by 7858 Related articles All 27 versions Web of Science: 1965 ☺

Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks

A Graves, S Fernández, F Gomez, J Schmidhuber

Proceedings of the 23rd international conference on Machine learning, 2006 - dl.acm.org

Many real-world sequence learning tasks require the prediction of sequences of labels from noisy, unsegmented input data. In speech recognition, for example, an acoustic signal is transcribed into words or sub-word units. Recurrent neural networks (RNNs) are powerful sequence learners that would seem well suited to such tasks. However, because they require pre-segmented training data, and post-processing to transform their outputs into label sequences, their applicability has so far been limited. This paper presents a

SHOW MORE ▾

☆ Save 99 Cite Cited by 6329 Related articles All 26 versions ☺

13

13

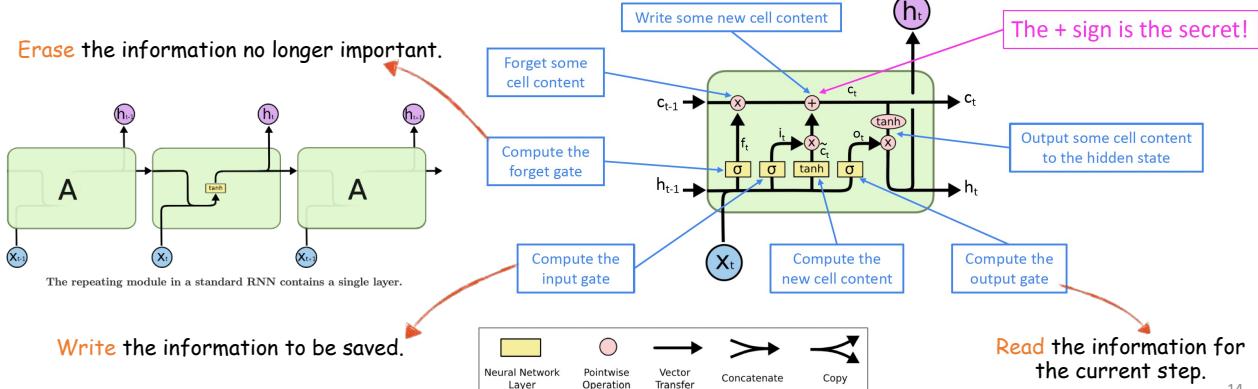
# LSTM Model Details

Reference: Stanford CS224N, Lecture 6:

<https://web.stanford.edu/class/cs224n/slides/cs224n-2024-lecture06-fancy-rnn.pdf>

<https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

- LSTM has a hidden state  $h^{(t)}$  and a cell state  $c^{(t)}$ , which stores **long-term information**.
  - The LSTM model can **read**, **erase**, and **write** information from the cell state, much like RAM.
  - Read/erase/write is controlled by three corresponding **gates**, which take values between 0 (**closed**) and 1 (**open**), and are **dynamically computed** based on the **current context**.



14

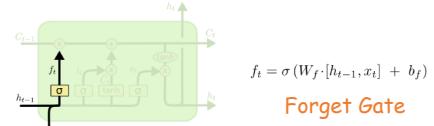
## LSTM Model Details

Reference: Stanford CS224N, Lecture 6:

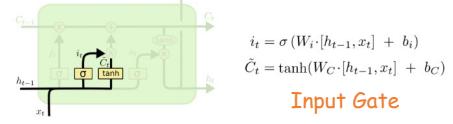
<https://web.stanford.edu/class/cs224n/slides/cs224n-2024-lecture06-fancy-rnn.pdf>

<https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

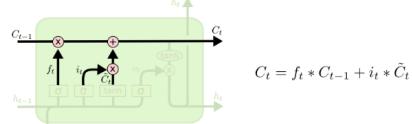
- Step 1: Decide how to forget the prior information.



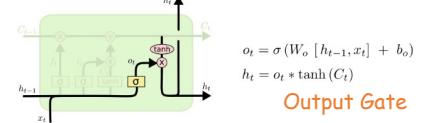
- Step 2: Decide the information needed to be stored to the cell.



- Step 3: Update the cell information.



- Step 4: Decide the output.



15

15

## LSTM for Sentiment Classification

		text	stars	sentiment
1411425	I went here and ordered the barbecue burger and I got the meat in medium well I had to admit it was really good although the fries got in the way it was salty and I had to dip the fries with mayo and ketchup all and all I would give the fries 35 out of 5 the service of the burger was right on time I like the look of the restaurant with the fun and relaxed look of the burger joint worth the visit especially if you just want to try a good burger		3	neg
971153	After calling 3 different companies number one were the only one that had the least wait time for our broken air in the middle of summer we called on monday morning and they were able to send a technician out that night at 6 pm when the technician came he determined it was our blower motor that was the problem and told us they will have to order the parts for it it will take couple days the repair will probably get done on wednesday and they will call us on tuesday to let us know the time they told us the repair will be between 11 am 2 pm on wednesday the repair guy pulled up in uhaul and after going up to the house to check the blower he told us he needed to go get more parts for it and come back in the afternoon around 5 pm still no sign of him so we called them back and asked them to come back again he came back and said he would be back soon come 7 pm still no sign so I called again and they told me they are on their way the guys did not show up until 845 pm took them about 30 minutes for the repair \n\nIn the end yes my ac was repaired in a timely manner I understand summer is the busiest time for them my issue is the lack of communication we were waiting and waiting and not knowing whats going on and had to keep calling them back also we never got a receipt the repair which they said they would email us		3	neg
1270461	went here for dinner and left very full and satisfied the family that runs the restaurant is originally from globe so the mexican food here is similar to what you would find in globe the food is delicious this restaurant has lots of options on the menu other than mexican food such as serrano and positas feels really atmospheric serving homemade no fried calamari i had the special tonight 2 chicken enchiladas with green sauce and rice beans the tortilla was perfect the enchiladas looked a bit mangled since chicken pieces were sticking out of the tortilla and the tortilla itself looked pressed and broken the enchiladas themselves were quite tasty with the green sauce and the rice and beans im usually a lightweight when it comes to finishing meals but in this case i cleaned my plate the salsa tastes good but is very watery which makes it hard to eat with chips we ordered iced teas and they were refilled promptly as needed id definitely be interested in going here again for some tasty and filling meals		4	pos

```

data['sentiment'] = ['pos' if (x>3) else 'neg' for x in data['stars']]
data['text'] = data['text'].apply((lambda x: re.sub('[^a-zA-Z0-9\s]', '', x)))
for idx, row in data.iterrows():
    row[0] = row[0].replace('rt', '')
data['text'] = [x.encode('ascii') for x in data['text']]

tokenizer = Tokenizer(nb_words=2500, lower=True, split=' ')
tokenizer.fit_on_texts(data['text'].values)
# print(tokenizer.word_index) # To see the dicstionary
X = tokenizer.texts_to_sequences(data['text'].values)
X = pad_sequences(X)
embed_dim = 128
lstm_out = 200
batch_size = 32

model = Sequential()
model.add(Embedding(2500, embed_dim, input_length = X.shape[1], dropout = 0.2))
model.add(LSTM(lstm_out, dropout_U = 0.2, dropout_W = 0.2))
model.add(Dense(2, activation='softmax'))
model.compile(loss = 'categorical_crossentropy', optimizer='adam', metrics = ['accuracy'])
print(model.summary())

```

16

16

## Evaluating Language Models

Reference: Stanford CS224N, Lecture 6:

<https://web.stanford.edu/class/cs224n/slides/cs224n-2024-lecture06-fancy-rnn.pdf>

<https://engineering.fb.com/2016/10/25/ml-applications/building-an-efficient-neural-language-model-over-a-billion-words/>

- **Perplexity:** The standard metric for evaluating a language model.

$$\text{perplexity} = \prod_{t=1}^T \left( \frac{1}{P_{\text{LM}}(\mathbf{x}^{(t+1)} | \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)})} \right)^{1/T}$$

Inverse probability of corpus, according to Language Model

Normalized by number of words

- Perplexity is the exponential of the cross-entropy loss, so the lower the better:

$$\text{Perplexity} = \prod_{t=1}^T \left( \frac{1}{\hat{y}_{\mathbf{x}^{(t+1)}}^{(t)}} \right)^{1/T} = \exp \left( \frac{1}{T} \sum_{t=1}^T -\log \hat{y}_{\mathbf{x}^{(t+1)}}^{(t)} \right) = \exp(J(\theta))$$

RNNs greatly improve perplexity over n-grams.

5-Gram  
DNN-based Methods  
(RNN, LSTM, etc.)

Model	Perplexity
Interpolated Kneser-Ney 5-gram (Chelba et al., 2013)	67.6
RNN-1024 + MaxEnt 9-gram (Chelba et al., 2013)	51.3
RNN-2048 + BlackOut sampling (Ji et al., 2015)	68.3
Sparse Non-negative Matrix factorization (Shazeer et al., 2015)	52.9
LSTM-2048 (Jozefowicz et al., 2016)	43.7
2-layer LSTM-8192 (Jozefowicz et al., 2016)	30
Ours small (LSTM-2048)	43.9
Ours large (2-layer LSTM-2048)	39.8

Table 2. Comparison on 1B word in perplexity (lower the better). Note that Jozefowicz et al., uses 32 GPUs for training. We only use 1 GPU.

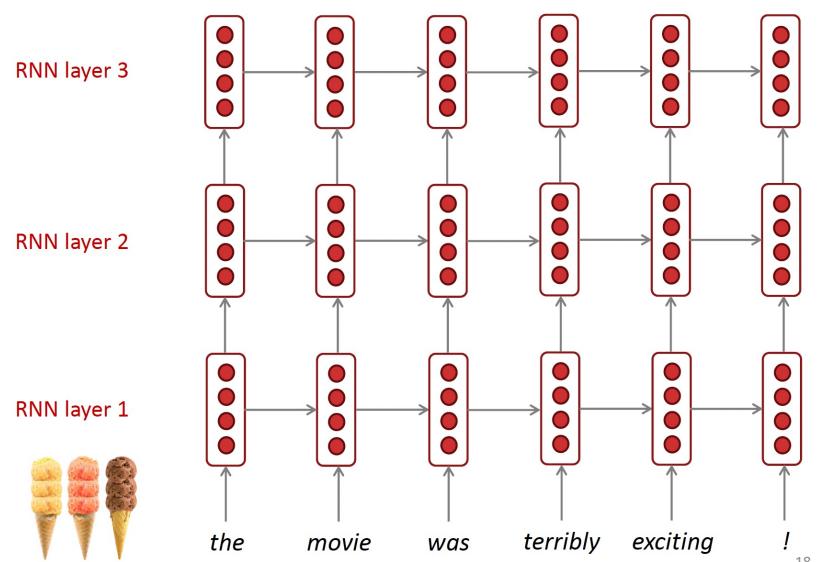
17

17

## Multi-layer (Stacked) RNN

Reference: Stanford CS224N, Lecture 6: <https://web.stanford.edu/class/cs224n/slides/cs224n-2024-lecture06-fancy-rnn.pdf>

- RNNs are already deep on time dimension.
- We can also deepen them in another dimension by applying multiple RNNs, allowing for more complex representations and achieving better performances.
- In neural machine translation, 2-4 layers is best for encoder RNNs; 4 layers is best for decoder RNNs.
- 2 layers is much better than 1, but 3 may be a little better than 2.
- Transformer-based nets are usually much deeper.



18

18

# Application of RNN: Inventory Management



https://pubsonline.informs.org/journal/mnsc

MANAGEMENT SCIENCE  
Vol. 69, No. 2, February 2023, pp. 759-773  
ISSN 0025-1909 (print), ISSN 1526-5501 (online)

## A Practical End-to-End Inventory Management Model with Deep Learning

Meng Qi,<sup>a</sup> Yuanxuan Shi,<sup>b</sup> Yongzhi Qi,<sup>c</sup> Chenxin Ma,<sup>d</sup> Rong Yuan,<sup>d</sup> Di Wu,<sup>d</sup> Zuo-Jun (Max) Shen<sup>a,\*</sup>

<sup>a</sup>SC Johnson College of Business, Cornell University, Ithaca, New York 14853; <sup>b</sup>Department of Electrical and Computer Engineering, University of California-San Diego, San Diego, California 92161; <sup>c</sup>IDeCorn Supply Chain V, Mountain View, California 94043; <sup>d</sup>IDeCorn Valley Research Center, Mountain View, California 94043; <sup>e</sup>College of Engineering, University of California-Berkeley, Berkeley, California 94720; <sup>f</sup>Faculty of Engineering & Faculty of Business and Economics, University of Hong Kong, Pokfulam, Hong Kong

\*Corresponding author.

Contact: mengd@cornell.edu (M. Qi); yyyshiheng@ucsd.edu (Y. Shi); qiyongzhi10@jhd.com (Y.Qi); chenxinma@ucsd.edu (C. Ma); rongyuan@jhd.com (R.Yuan); di.wu@jhd.com (D.Wu); shenzh@berkeley.edu (Z.J. Shen)

https://doi.org/10.1287/mnsc.2022.4564

Copyright © 2022 INFORMS

**Abstract:** We investigate a data-driven multiperiod inventory replenishment problem with uncertain demand and vendor lead time (VLT) with accessibility to a large quantity of historical data. Different from the traditional two-step predict-then-optimize (PTO) solution framework, we propose a one-step end-to-end (E2E) framework that uses deep learning models to output the suggested replenishment amount directly from input features without any intermediate step. The E2E model is trained to capture the behavior of the optimal dynamic replenishment scheme under historical observations without any prior assumptions on the distribution of the data and does not consider any specific domain knowledge. Numerical experiments using real data from one of the leading e-commerce companies, we demonstrate the advantages of the proposed E2E model over conventional PTO frameworks. We also conduct a field experiment with IDeCorn, and the results show that the new algorithm significantly reduces the total cost of holding inventories and improves rates substantially compared with ID's current practice. For the supply chain and management industry, our E2E model shortens the decision process and provides an automatic inventory management service with the power to generalize and scale. The concept of E2E, which makes the best information available for a certain goal, can also be used in practice for other supply chain management circumstances.

**History:** Accepted by Hanan Naarazadeh, big data analytics.

**Funding:** This research was supported by the National Key Research and Development Program of China [Grant 2018YFB170000] and National Science Foundation of China [Grants 71991462 and 71971310].

**Supplemental Material:** The online data are available at https://doi.org/10.1287/mnsc.2022.4564.

**Keywords:** end-to-end decision-making • inventory management • deep learning • e-commerce

- RNN is not that frequently used in business research, because it does NOT directly produce text representations.

- Use multi-quantile RNNs to provide end-to-end predictions from features to the optimal inventory decisions, whereas most of the literature applies the predict-then-optimize paradigm.
- A field experiment shows that the e2e approach substantially reduces the inventory costs compared with some naïve benchmarks.

## A practical end-to-end inventory management model with deep learning

M.Qi, Y.Shi, Y.Qi, C.Ma, R.Yuan, D.Wu, Z.J.Shen

Management Science, 2023 pubsonline.informs.org

We investigate a data-driven multiperiod inventory replenishment problem with uncertain demand and vendor lead time (VLT) with accessibility to a large quantity of historical data. Different from the traditional two-step predict-then-optimize (PTO) solution framework, we propose a one-step end-to-end (E2E) framework that uses deep learning models to output the suggested replenishment amount directly from input features without any intermediate step. The E2E model is trained to capture the behavior of the optimal dynamic replenishment scheme under historical observations without any prior assumptions on the distribution of the data and does not consider any specific domain knowledge.

SHOW MORE ▾

☆ Save 99 Cite Cited by 65 Related articles All 4 versions Web of Science: 5

19

19

# Application of RNN: Detecting FTD



Psychiatry Research 304 (2021) 114135

Contents lists available at ScienceDirect

Psychiatry Research

journal homepage: www.elsevier.com/locate/psychres



Detecting formal thought disorder by deep contextualized word representations

Justyna Sarzyńska-Wawer<sup>1,a</sup>, Aleksander Wawer<sup>1,b</sup>, Aleksandra Pawlak<sup>2,c</sup>,  
Julia Szmaynowska<sup>2,d</sup>, Izabela Stefanik<sup>1</sup>, Michał Jarkiewicz<sup>3,e,f</sup>, Lukasz Okruszek<sup>1</sup>

<sup>1</sup> Institute of Psychology, Polish Academy of Sciences, Warsaw, Poland

<sup>2</sup> Institute of Computer Science, Polish Academy of Sciences, Warsaw, Poland

<sup>3</sup> University of Social Sciences and Humanities, Warsaw, Poland

<sup>a,b</sup> Institute of Psychiatry and Neurology, Warsaw, Poland

ARTICLE INFO

ABSTRACT

Computational linguistics has enabled the introduction of objective tools that measure some of the symptoms of schizophrenia, including the coherence of speech associated with formal thought disorder (FTD). Our goal was to investigate whether neural network based utterance embeddings are more accurate in detecting FTD than models based on individual indicators. The present research used a comprehensive Embeddings from Language Models (ELMo) approach to represent interviews with patients suffering from schizophrenia (N=35) and with healthy people (N=35). We compare its results to the approach described by Bedi et al. (2015), referred to here as the coherence model. Evaluation was also performed using the rating scale for the Assessment of Formal Thought Disorders (TLC) and Language and Communication (TLG) questions. The ELMo model achieved an accuracy of 80% in distinguishing patients from healthy people. Previously used coherence models were less accurate at 70%. The classifying clinician was accurate 74% of the time. Our analysis shows that both ELMo and TLC are sensitive to the symptoms of disorganization in patients. In this study methods using text representations from language models were more accurate than those based solely on the assessment of FTD, and can be used as measures of disordered language that complement human clinical ratings.

- How to detect formal thought disorder (FTD)?
- Embeddings from LSTM language models (ELMo) can more accurately detect/predict FTD than individual indicators (benchmark: coherence model, which can somehow be viewed as traditional NLP method).
- Accuracy (N=70, 35 healthy and 35 patients):
  - ELMo: 80%
  - Coherence models: <70%
  - Clinician: 74%

## [HTML] Detecting formal thought disorder by deep contextualized word representations

J.Sarzyńska-Wawer, A.Wawer, A.Pawlak... - Psychiatry ..., 2021 - Elsevier

Computational linguistics has enabled the introduction of objective tools that measure some of the symptoms of schizophrenia, including the coherence of speech associated with formal thought disorder (FTD). Our goal was to investigate whether neural network based utterance embeddings are more accurate in detecting FTD than models based on individual indicators. The present research used a comprehensive Embeddings from Language Models (ELMo) approach to represent interviews with patients suffering from schizophrenia ...

☆ Save 99 Cite Cited by 14561 Related articles All 24 versions Web of Science: 74

20

20

10

## Agenda

- Vanilla Recurrent Neural Nets (RNN)
- Long Short-Term Memory (LSTM)
- Sequence-to-sequence (Seq2seq)

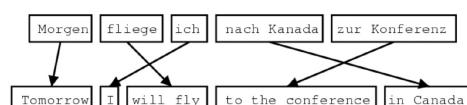
21

21

## Neural Machine Translation (NMT)

Reference: Stanford CS224N, Lecture 6: <https://web.stanford.edu/class/cs224n/slides/cs224n-2024-lecture06-fancy-rnn.pdf>

- NMT is a way to do machine translation with a single end-to-end neural network: Sequence-to-sequence (**seq2seq**), which involves **2 RNNs**.
- Machine translation is **highly nontrivial** and once was a huge research field in CS and NLP.



1519年600名西班牙人在墨西哥登陆，去征服几百万人口的阿兹特克帝国，初次交锋他们损兵三分之二。

In 1519, six hundred Spaniards landed in Mexico to conquer the Aztec Empire with a population of a few million. They lost two thirds of their soldiers in the first clash.

[translate.google.com \(2009\)](#): 1519 600 Spaniards landed in Mexico, millions of people to conquer the Aztec empire, the first two-thirds of soldiers against their loss.

[translate.google.com \(2013\)](#): 1519 600 Spaniards landed in Mexico to conquer the Aztec empire, hundreds of millions of people, the initial confrontation loss of soldiers two-thirds.

[translate.google.com \(2015\)](#): 1519 600 Spaniards landed in Mexico, millions of people to conquer the Aztec empire, the first two-thirds of the loss of soldiers they clash.

22

22

## Seq2Seq for NMT

Reference: Stanford CS224N, Lecture 6: <https://web.stanford.edu/class/cs224n/slides/cs224n-2024-lecture06-fancy-rnn.pdf>

- Seq2seq is a Conditional Language Model:

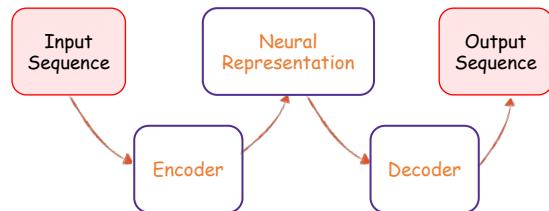
- Predicting the next word of the target sentence  $y$  conditioned on the source sentence  $x$  and prior texts.

$$P(y|x) = P(y_1|x) P(y_2|y_1, x) P(y_3|y_1, y_2, x) \dots P(y_T|y_1, \dots, y_{T-1}, x)$$

Probability of next target word, given target words so far and source sentence  $x$

- Encoder-decoder architecture: Encoder takes input and produces a neural representation; Decoder produces output based on that neural representation.

- Seq2seq: both input and output are sequences.
- Summarization: Long text  $\rightarrow$  short text
- Dialogue: previous utterances  $\rightarrow$  next utterance
- Parsing: Input text  $\rightarrow$  output parse as a sequence
- Code generation  $\rightarrow$  Natural language  $\rightarrow$  Python code

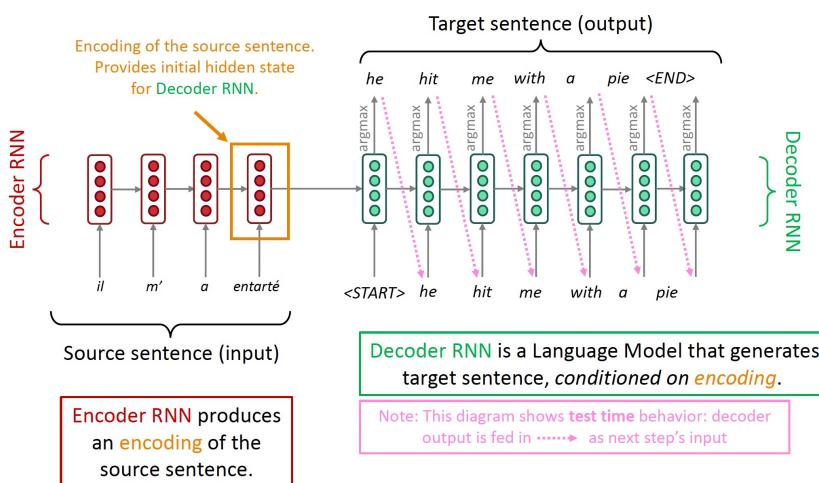


23

23

## Seq2Seq Architecture

Reference: Stanford CS224N, Lecture 6: <https://web.stanford.edu/class/cs224n/slides/cs224n-2024-lecture06-fancy-rnn.pdf>



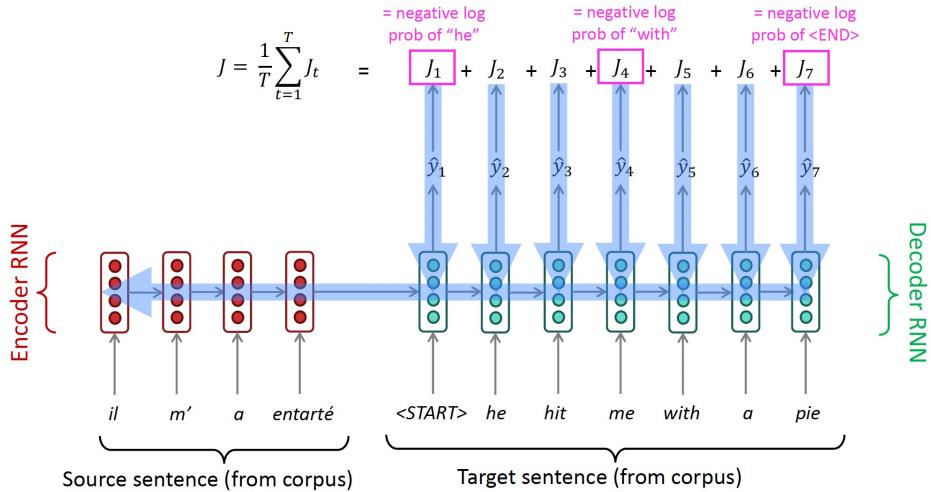
[Sequence to sequence learning with neural networks](#)  
I.Sutskever, O.Vinyals, Q.V.Le - Advances in neural ... 2014 - proceedings.neurips.cc  
Abstract Deep Neural Networks (DNNs) are powerful models that have achieved excellent performance on difficult learning tasks. Although DNNs work well whenever large labeled training sets are available, they cannot be used to map sequences to sequences. In this paper, we present a general end-to-end approach to sequence learning that makes minimal assumptions on the sequence structure. Our method uses a multilayered Long Short-Term Memory (LSTM) to map the input sequence to a vector of a fixed dimensionality, and then ...  
☆ Save 59 Cite Cited by 25043 Related articles All 28 versions

24

24

## Seq2Seq Training

Reference: Stanford CS224N, Lecture 6: <https://web.stanford.edu/class/cs224n/slides/cs224n-2024-lecture06-fancy-rnn.pdf>



Seq2seq is optimized as a single system. Backpropagation operates “end-to-end”.

25

25

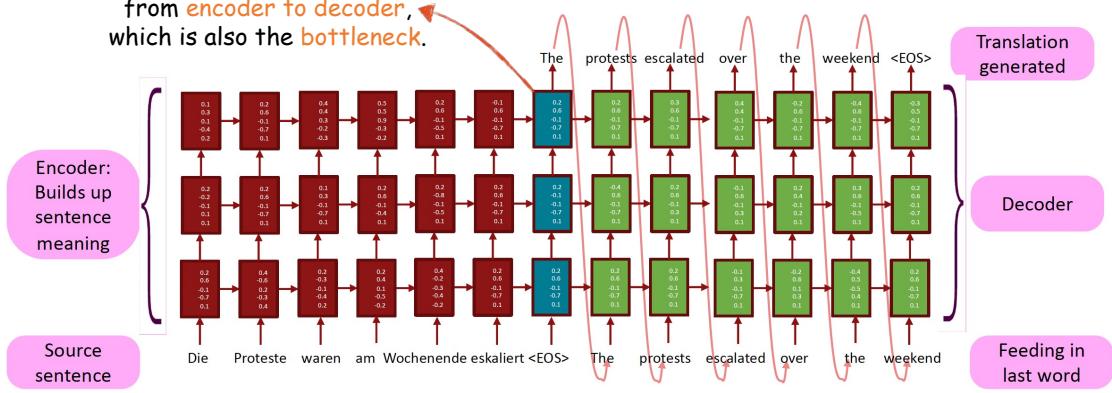
## Multi-Layer Seq2Seq

Reference: Stanford CS224N, Lecture 6: <https://web.stanford.edu/class/cs224n/slides/cs224n-2024-lecture06-fancy-rnn.pdf>

[Sutskever et al. 2014; Luong et al. 2015]

The hidden states from RNN layer  $i$  are the inputs to RNN layer  $i+1$

Conditioning: Information flow from encoder to decoder, which is also the bottleneck.



26

26