

Sequence Models Week 1

Examples of sequence data

Speech recognition		→	"The quick brown fox jumped over the lazy dog."
Music generation		→	
Sentiment classification	"There is nothing to like in this movie."	→	★☆☆☆☆
DNA sequence analysis	→ AGCCCCCTGTGAGGAACCTAG	→	AGCCCCCTGTGAGGAACCTAG
Machine translation	Voulez-vous chanter avec moi?	→	Do you want to sing with me?
Video activity recognition		→	Running
Name entity recognition	→ Yesterday, Harry Potter met Hermione Granger.	→	Yesterday, Harry Potter met Hermione Granger.

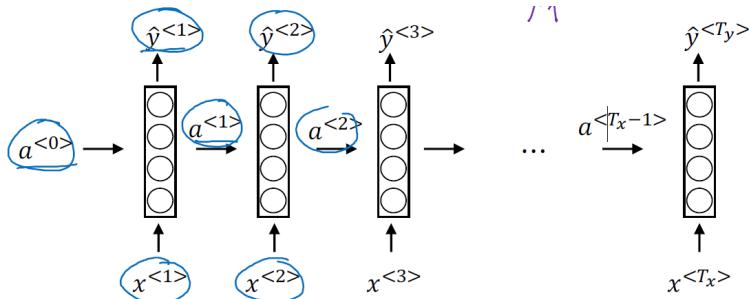
Recurrent Neural Network Model

why not a standard network?

problem:

- Inputs, outputs can be different lengths in different examples
- Doesn't share features learned across different positions of text.

RNN model:



Drawback: 1. 用前向的模型去预测后续的。
2. $x^{<1>} \rightarrow x^{<2>} \rightarrow \dots \rightarrow x^{<T_x>} \rightarrow \dots$ 只能使用 $x^{<1>} \rightarrow x^{<2>} \rightarrow \dots \rightarrow x^{<T_x>} \rightarrow \dots$ 的信息，但前面的 $x^{<1>} \rightarrow x^{<2>} \rightarrow \dots \rightarrow x^{<T_x>} \rightarrow \dots$ 的信息也是有用的。

forward propagation:

$$\begin{aligned}
 a^{<0>} &= \vec{0} \quad | \quad a^{<t>} = g_1(W_{aa} a^{<t-1>} + W_{ax} x^{<t>} + b_a) \leftarrow \tanh / \text{relu} \\
 \hat{y}^{<t>} &= g_2(W_{ya} a^{<t>} + b_y) \leftarrow \text{sigmoid} / \text{softmax} \\
 a^{<t>} &= g_t(W_{aa} a^{<t-1>} + W_{ax} x^{<t>} + b_a) \\
 \hat{y}^{<t>} &= g_t(W_{ya} a^{<t>} + b_y)
 \end{aligned}$$

Try to simplify it to suit the more complex network

$$a^{<t>} = g_t(W_a a^{<t-1>} + W_x x^{<t>} + b_a) \rightarrow a^{<t>} = g(W_a [a^{<t-1>}, x^{<t>}] + b_a)$$

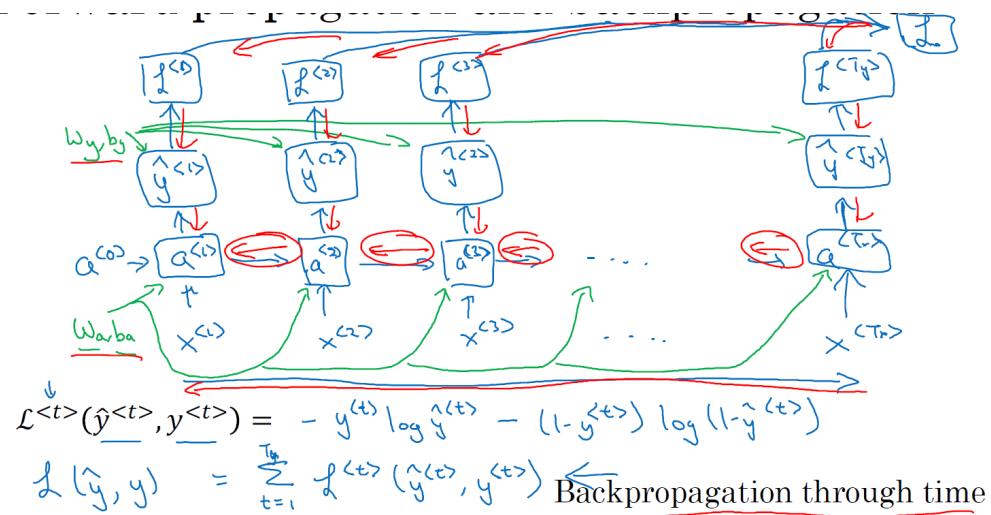
$$g^{<t>} = g(W_y a^{<t>} + b_y)$$

$$\hookrightarrow \hat{y}^{<t>} = g(W_y a^{<t>} + b_y)$$

$$W_a = [W_{aa}; W_{ax}]$$

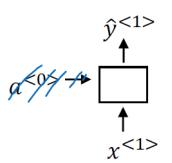
$$[a^{<t-1>}, x^{<t>}] = \begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix}$$

Backpropagation Through Time



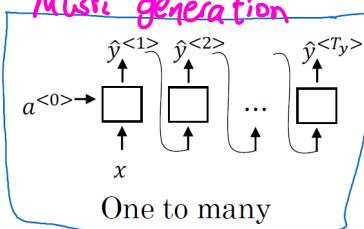
Different Types of RNNs

Summary of RNN types



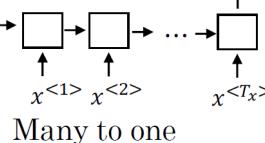
One to one

Music generation

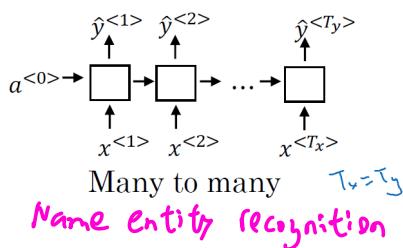


One to many

Sentiment classification

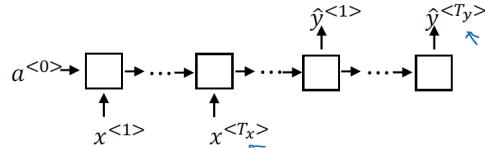


Many to one



Many to many

Name entity recognition

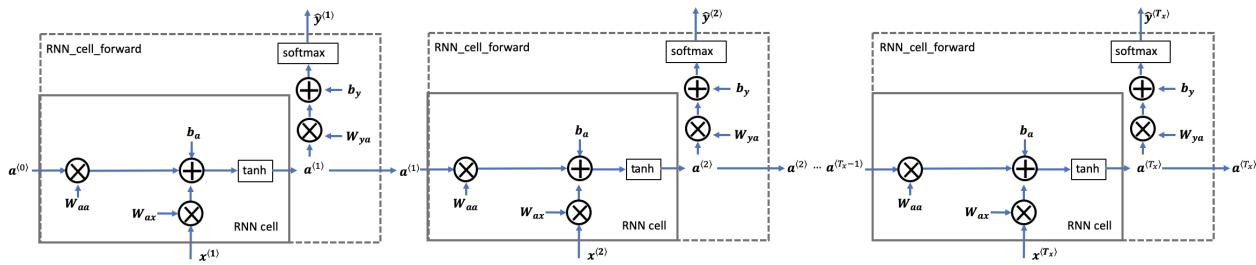
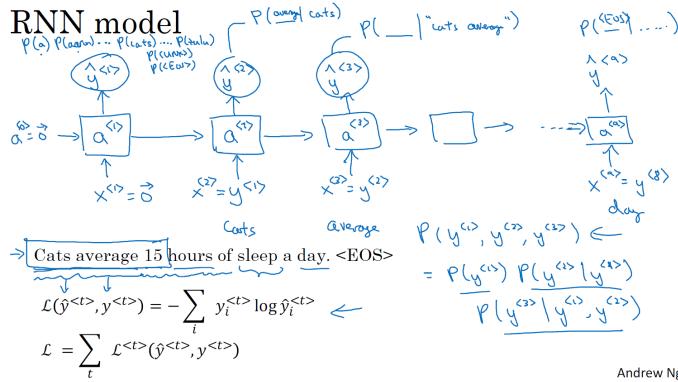


Many to many

Machine translation

language model and Sequence Generation

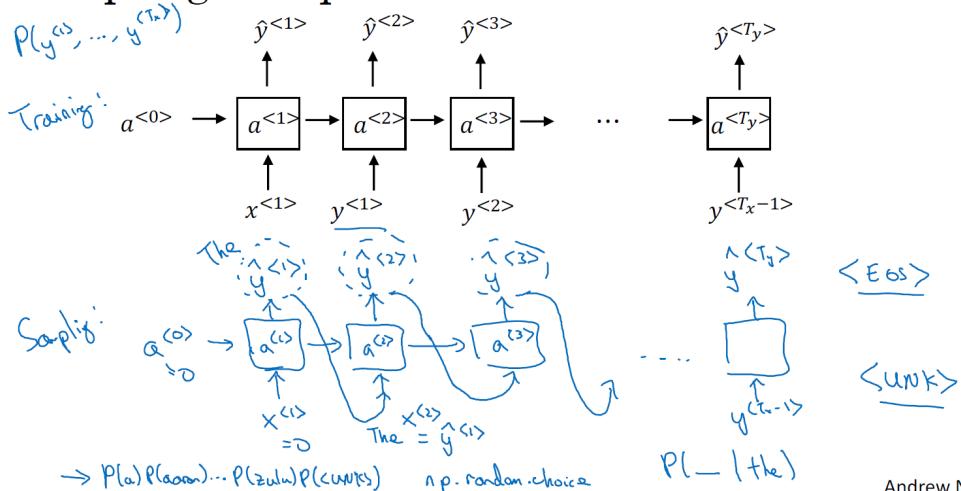
Training set: large corpus of English text.



Sampling Novel Sequences

Character-level
word-level

Sampling a sequence from a trained RNN



Vanishing Gradient with RNNs

Explode / Vanish

The basic RNNs cannot catch the long-range dependency.

Why?

Recap: Why vanishing gradient?

① 梯度消失和梯度爆炸都可能导致梯度消失，因为梯度爆炸时梯度的值会很大。

导致梯度 NP 小 (大)。

② Sigmoid 是函数 $f(x) \in [0, 1]$

③ 和前层梯度 Δ_t 相加 \rightarrow $\begin{cases} \text{ReLU} & \text{if sigmoid} \\ \text{BN} & \text{if ResNet} \end{cases}$

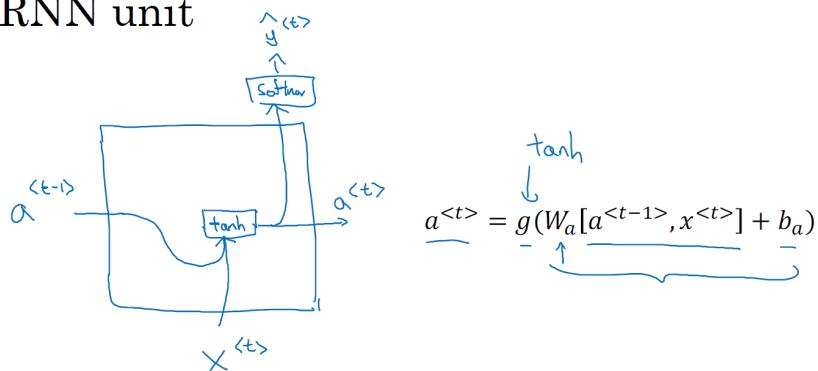
How to solve?

Explode \rightarrow gradient clipping

Vanish \rightarrow ? GRU / LSTM

Gated Recurrent Unit (GRU)

RNN unit



GRU *

C = memory cell

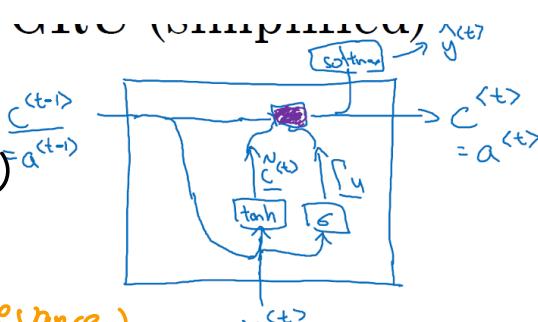
$$c^{<t>} = a^{<t>}$$

$$\rightarrow \tilde{c}^{<t>} = \tanh(W_k[\Gamma_r * c^{t-1}, x^{t-1}] + b_r)$$

$$\rightarrow \Gamma_u = \sigma(W_u[c^{t-1}, x^{t-1}] + b_u)$$

$$\rightarrow \Gamma_r = \sigma(W_r[c^{t-1}, x^{t-1}] + b_r) \quad (\text{relevance}).$$

$$C^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * C^{t-1}$$



Long Short Term Memory (LSTM)

GRU

$$\tilde{c}^{<t>} = \tanh(W_c[\Gamma_r * c^{<t-1>}, x^{<t>}] + b_c)$$

$$\Gamma_u = \sigma(W_u[c^{<t-1>}, x^{<t>}] + b_u)$$

$$\Gamma_r = \sigma(W_r[c^{<t-1>}, x^{<t>}] + b_r)$$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>}$$

$$a^{<t>} = c^{<t>}$$

LSTM

$$\tilde{c}^{<t>} = \tanh(W_c[a^{<t-1>}, x^{<t>}] + b_c)$$

$$\Gamma_u = \sigma(W_u[a^{<t-1>}, x^{<t>}] + b_u)$$

$$\Gamma_f = \sigma(W_f[a^{<t-1>}, x^{<t>}] + b_f)$$

forget gate

$$\Gamma_o = \sigma(W_o[a^{<t-1>}, x^{<t>}] + b_o)$$

update gate

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + \Gamma_f * c^{<t-1>}$$

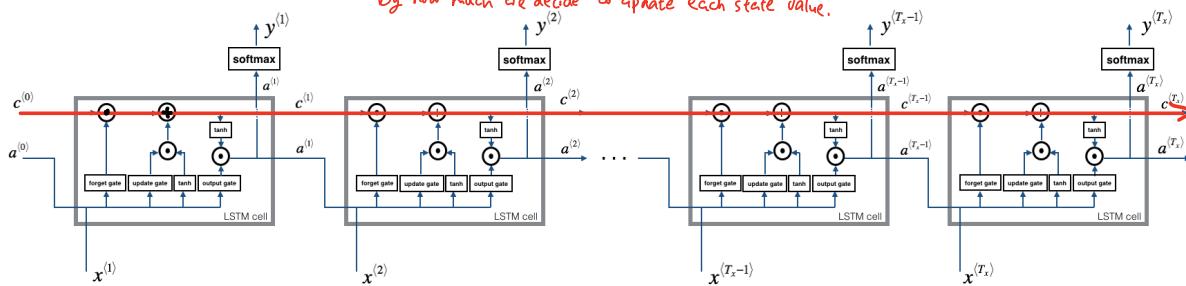
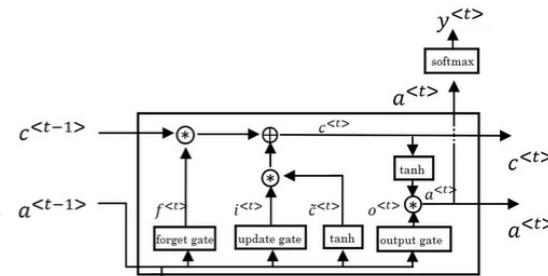
output gate

$$a^{<t>} = \Gamma_o * \tanh(c^{<t>})$$

$$\begin{cases} \tilde{c}^{<t>} = \tanh(W_c[a^{<t-1>}, x^{<t>}] + b_c) \\ \Gamma_u = \sigma(W_u[a^{<t-1>}, x^{<t>}] + b_u) \\ \Gamma_f = \sigma(W_f[a^{<t-1>}, x^{<t>}] + b_f) \\ \Gamma_o = \sigma(W_o[a^{<t-1>}, x^{<t>}] + b_o) \\ c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + \Gamma_f * c^{<t-1>} \\ a^{<t>} = \Gamma_o * \tanh(c^{<t>}) \end{cases}$$

*multiply the old state by Γ_f , forgetting the things we decided to forget earlier. Then we add $\Gamma_u * \tilde{c}^{<t>}$. This is a new candidate values, scaled by how much we decide to update each state value.*

next



- a **forget gate**: decide which inputs units should be remembered and passed along, It is a tensor with values between 0 and 1
- an **update gate**: again a tensor containing values between 0 and 1. It decides on what information to throw away, and what new information to add.
- **output gate**: decide what gets as the output of the time step.

Bi-directional RNNs

Deep RNNs

Sequence Model Week 2 Word Representation

Word representation

Transfer learning and word embeddings

1. Learn word embeddings from large text corpus.
(Or download pre-trained embeddings online)
2. Transfer embedding to new task with smaller training set
3. Continue to fine tune the word embeddings with new data.

Properties of word embeddings

t-SNE

$$e_{man} - e_{woman} = e_{king} - e_{woman}$$

Find word $w \Rightarrow \arg \max_w \text{Similarity}(e_w, e_{king} + e_{woman} - e_{woman})$

cosine similarity

$$\text{Sim}(u, v) = \frac{u^T v}{\|u\|_2 \|v\|_2}$$

Embedding matrix

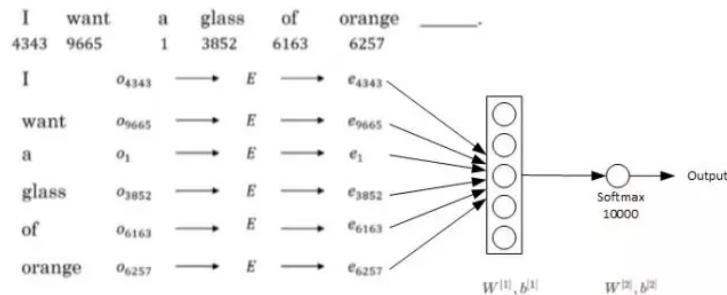
$$E \quad \begin{bmatrix} \text{orange} \\ \vdots \\ \text{orange} \end{bmatrix}_{300 \times 10000} \quad \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}_{10000 \times 1} \quad \text{One-hot} \quad e_{\text{orange}} \quad 300 \times 1$$

For a given word w , the embedding matrix E is a matrix that maps its 1-hot representation to its embedding e_w as follows:

$$e_w = E \cdot O_w$$

Learning Word Embeddings: Word2Vec & Glove

Learning Word Embeddings



Word2Vec

Word2Vec is a framework aimed at learning word embeddings by estimating the likelihood that a given word is surrounded by other words. Popular models include Skip-gram, Negative Sampling and CBOW

Skip-gram

The skip-gram Word2Vec model is a supervised learning task that learns word embeddings by accessing the likelihood of any given target word t happening with a context word c .

$$p(c|t) = \frac{\exp(\theta_t^T e_c)}{\sum_{j=1}^{|V|} \exp(\theta_j^T e_c)}$$

$$\text{Softmax: } \sigma(z)_i = \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}}$$

where θ_c is a parameter associated with c .

Drawback: summing over the whole vocabulary in the denominator of the softmax part makes this model computationally expensive

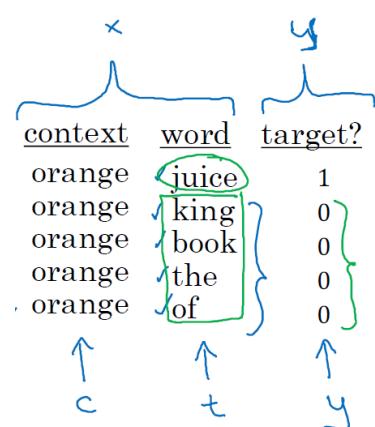
Some methods to improve: hierarchical softmax model.

Negative Sampling

It is set of binary classifiers using logistic regression that aim at assessing how a given context and a given target words are likely to appear simultaneously. with the models being trained on sets of K negative examples and 1 positive example

Example

$$p(y=1 | c, t) = \sigma(\theta_t^T e_c)$$



,

Sentiment classification

Sentiment classification problem

x —————→ y

The dessert is excellent.	★★★★★☆
Service was quite slow.	★★☆☆☆☆
Good for a quick meal, but nothing special.	★★★☆☆☆
Completely lacking in good taste, good service, and good ambience.	★☆☆☆☆☆