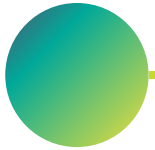




VietAI

Recurrent Neural Network



VietAI teaching team

1 Giới thiệu bài toán

This is good!

Not really bad

Awful



Negative



Neutral






Positive

1 Giới thiệu bài toán

ANH - ĐÃ PHÁT HIỆN ANH VIỆT PHÁP ▾





Today I learn how a recurrent neural network works ✕

  50/5000 

VIỆT ANH TRUNG (GIẢN THỂ) ▾

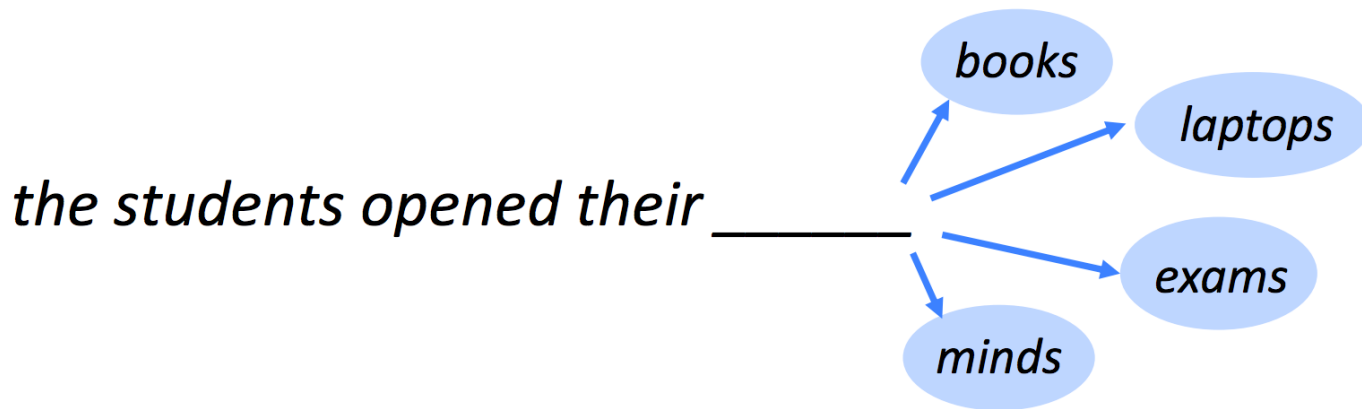
今天，我学习了循环神经网络的工作原理 ☆

Jīntiān, wǒ xuéxíle xúnhuán shénjīng wǎngluò de gōngzuò yuánlǐ

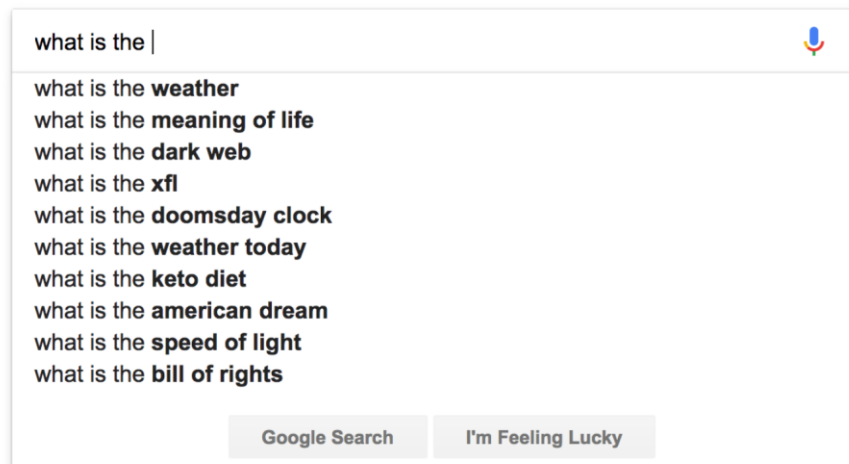
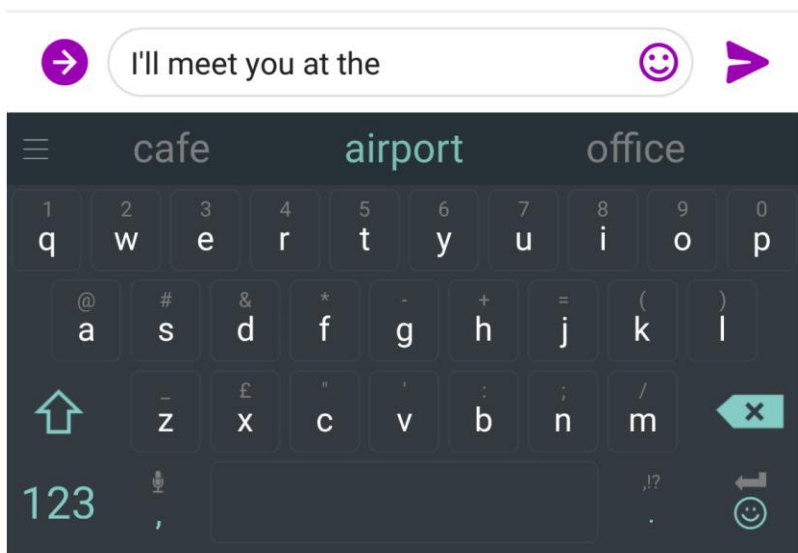
1 Giới thiệu bài toán

Language model

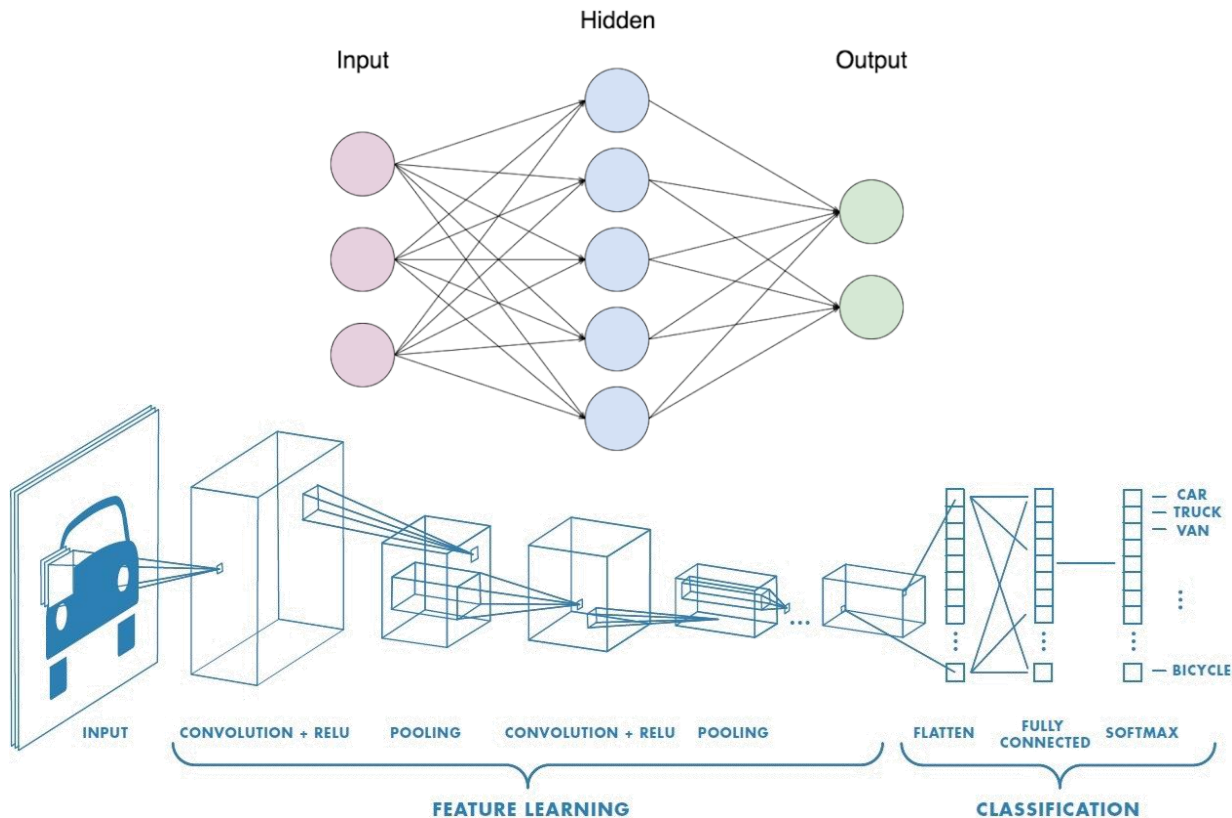


$$P(w_m | w_1, \dots, w_{m-1})$$

1 Giới thiệu bài toán



1 Giới thiệu bài toán



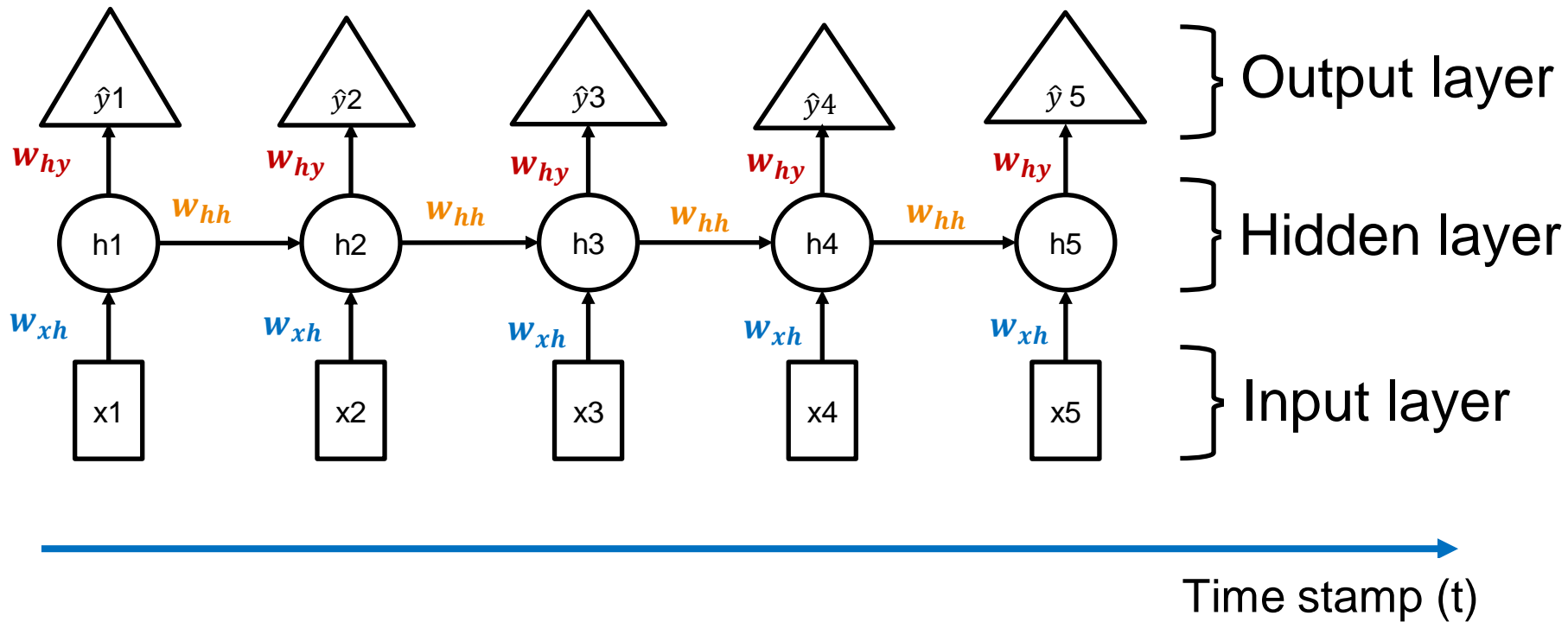
Nội dung

1. Cấu trúc Recurrent Neural Network (RNN)
2. Cách hoạt động
3. Backpropagation through time
4. Vanishing/Exploding gradients

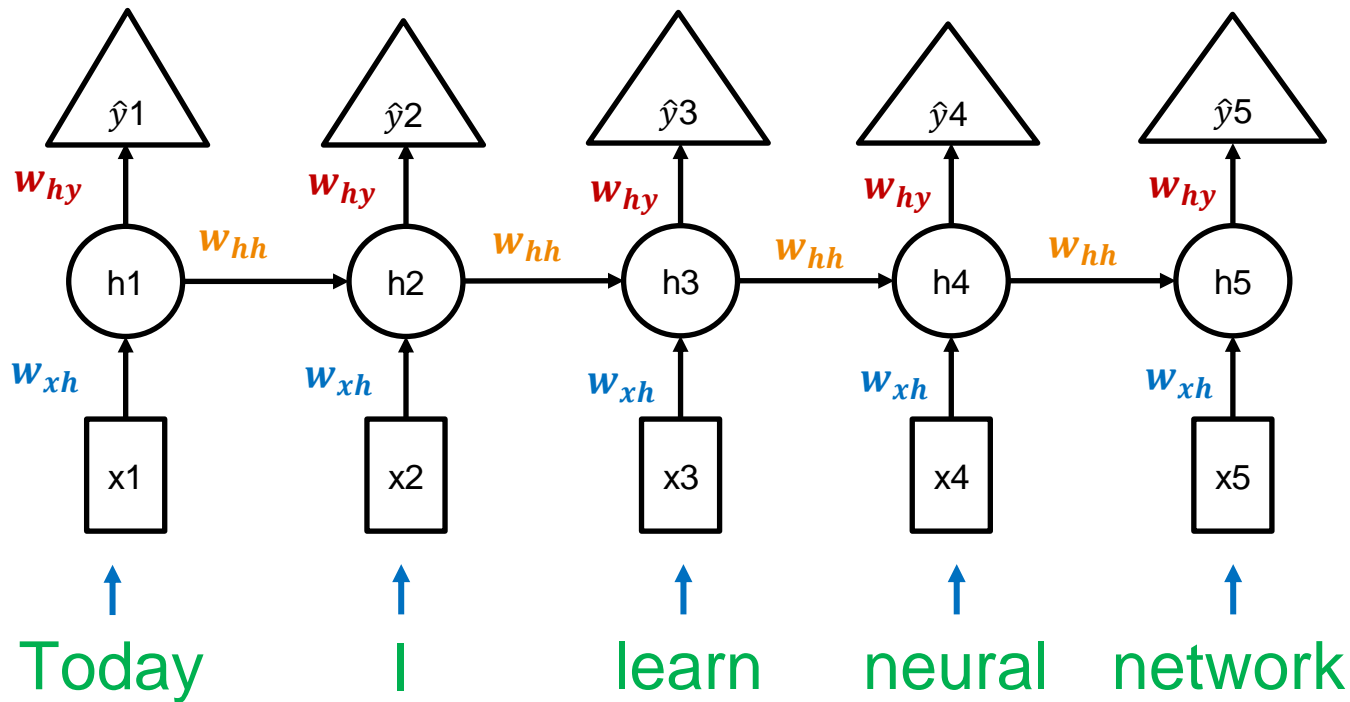
Nội dung

1. Cấu trúc Recurrent Neural Network (RNN)
2. Cách hoạt động
3. Backpropagation through time
4. Vanishing/Exploding gradients

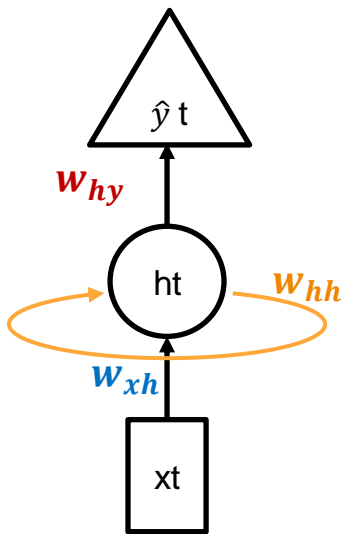
1 Cấu trúc RNN



1 Cấu trúc RNN



1 Cấu trúc RNN

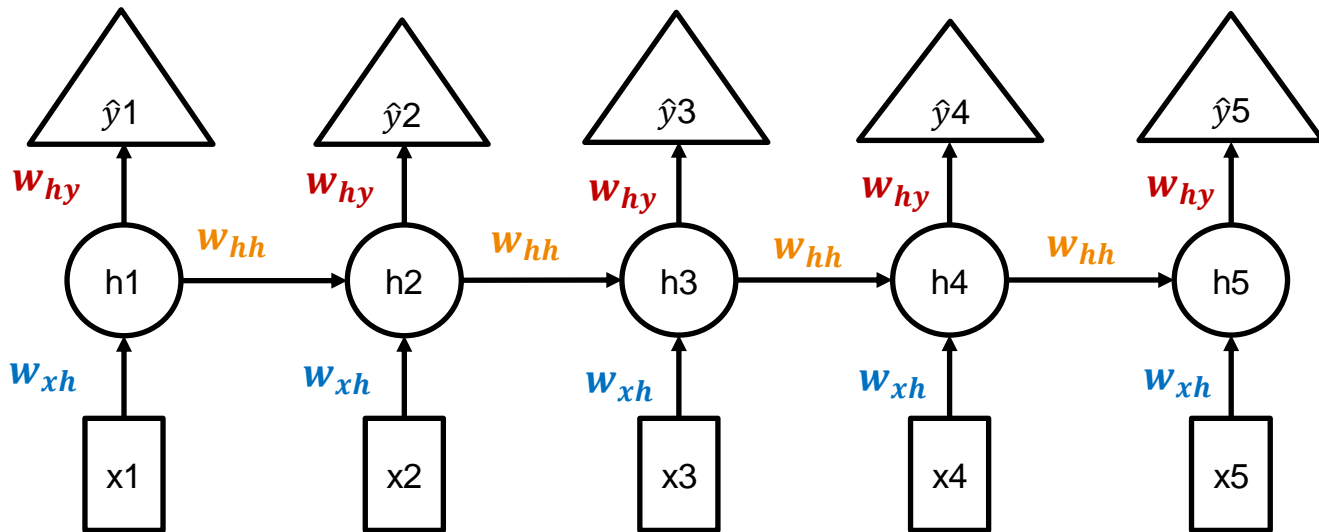


Today I learn neural network

Nội dung

1. Cấu trúc Recurrent Neural Network (RNN)
2. **Cách hoạt động**
3. Backpropagation through time
4. Vanishing/Exploding gradients

2 Cách hoạt động của RNN



Step 1: $h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t + b_h)$

Step 2: $\hat{y}_t = \text{softmax}(W_{hy}h_t + b_y)$

$$\tanh(x) = \frac{1 - e^{-2x}}{1 + e^{-2x}}$$

$$\text{softmax}(x_1) = \frac{e^{x_1}}{e^{x_1} + e^{x_2} + \dots + e^{x_n}}$$

2 Cách hoạt động của RNN

Ví dụ: Tìm giá trị \hat{y}_2 , biết:

$$x_1 = [0.1 \quad 0.2]$$

$$x_2 = y_1$$

$$w_{xh} = \begin{bmatrix} 0 & 1 & 0.1 \\ 0.3 & 0.5 & 1 \end{bmatrix}$$

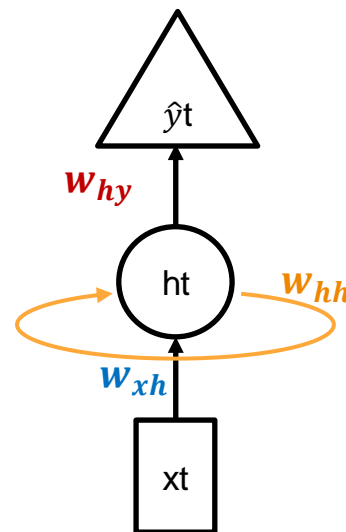
$$w_{hh} = \begin{bmatrix} 1 & 0.1 & 0.2 \\ 0.5 & 0.5 & 1 \\ 0 & 1 & 0 \end{bmatrix}$$

$$w_{hy} = \begin{bmatrix} 1 & 1 \\ 0.1 & 1 \\ 1 & 0.5 \end{bmatrix}$$

$$h_0 = [0 \quad 0 \quad 0]$$

$$b_h = [0.1 \quad 0.1 \quad 0]$$

$$b_y = [0 \quad 0.5]$$



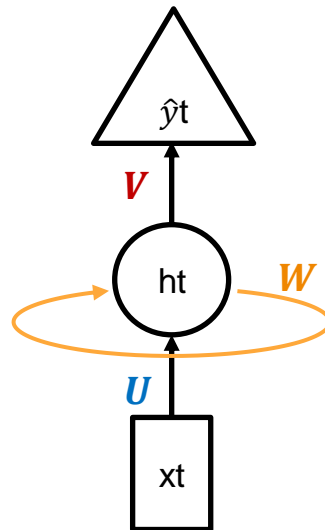
2

Cách hoạt động của RNN

Giải:

2 Quy định lại ký hiệu

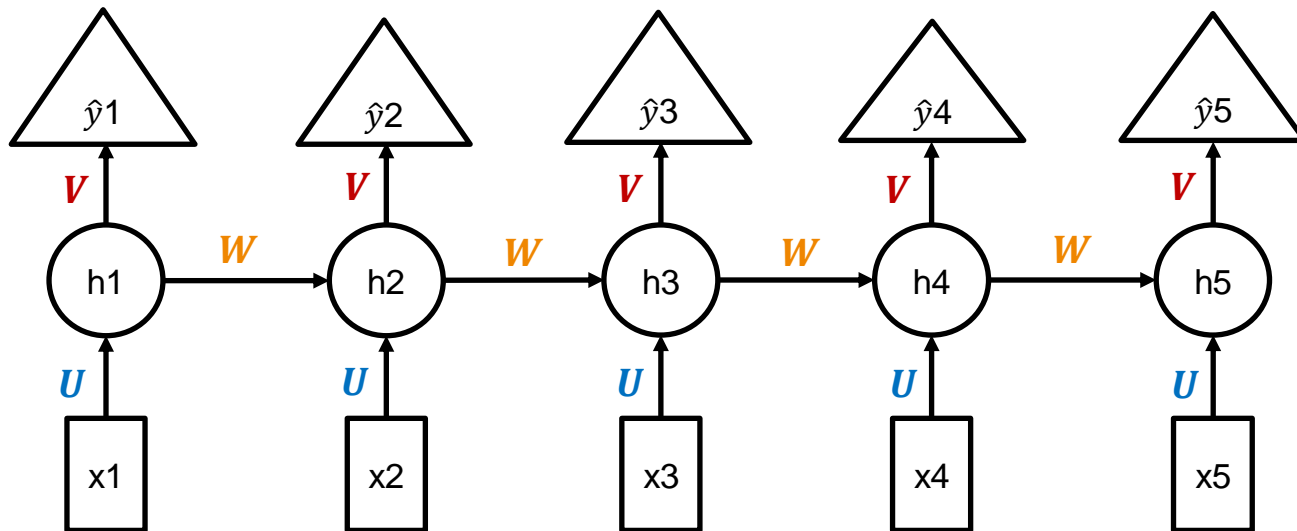
$$\begin{aligned} w_{xh} &\equiv U \\ w_{hh} &\equiv W \\ w_{hy} &\equiv V \end{aligned}$$



Nội dung

1. Cấu trúc Recurrent Neural Network (RNN)
2. Cách hoạt động
3. **Backpropagation through time**
4. Vanishing/Exploding gradients

3 Backpropagation through time



$$Loss^{<t>}(\hat{y}^{<t>}, y^{<t>}) = -y^{<t>} \log \hat{y}^{<t>} - (1 - y^{<t>}) \log(1 - \hat{y}^{<t>})$$

$$\mathbf{L}(\hat{y}, y) = \sum_{t=1}^T Loss^{<t>}(\hat{y}^{<t>}, y^{<t>})$$

3 Backpropagation through time

- Đạo hàm hàm Loss (tương tự Lecture 8)
- Tìm gradients tương ứng cho mỗi ma trận trọng số U , W , V

$$U = U - \alpha \Delta U$$

$$W = W - \alpha \Delta W$$

$$V = V - \alpha \Delta V$$

- Cần tìm $\frac{\partial L}{\partial V}$, $\frac{\partial L}{\partial W}$, $\frac{\partial L}{\partial U}$

3 Backpropagation through time

Đạo hàm riêng hàm L trên biến V

$$\frac{\partial L}{\partial V} = \sum_{t=1}^T \frac{\partial Loss^{<t>}}{\partial V}$$

Tại một thời điểm t nhất định:

$$= [-y \log \hat{y} - (1 - y) \log(1 - \hat{y})]'$$

$$= [-y \log(\text{softmax}(V h_t + b_y)) - (1 - y) \log(1 - \text{softmax}(V h_t + b_y))]'$$

→ Sử dụng đạo hàm hàm hợp

Step 1: $h_t = \tanh(W_{hh} h_{t-1} + W_{xh} x_t + b_h)$

Step 2: $\hat{y}_t = \text{softmax}(W_{hy} h_t + b_y)$

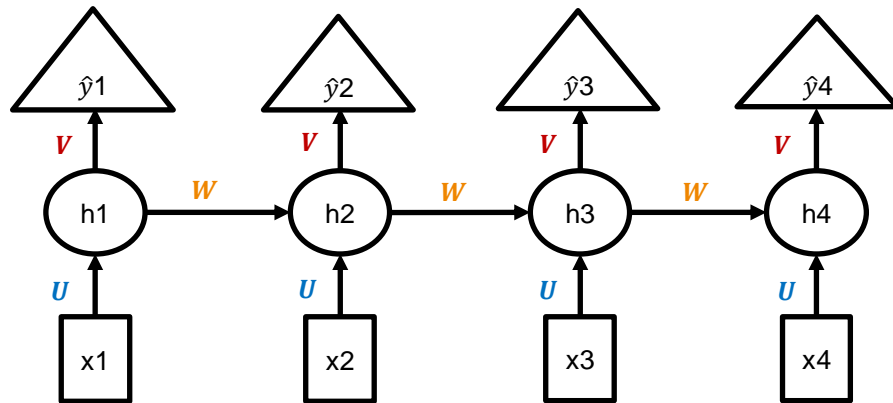
3 Backpropagation through time

Đạo hàm riêng hàm L trên biến W

$$\frac{\partial L}{\partial W} = \sum_{t=1}^T \frac{\partial Loss^{<t>}}{\partial W}$$

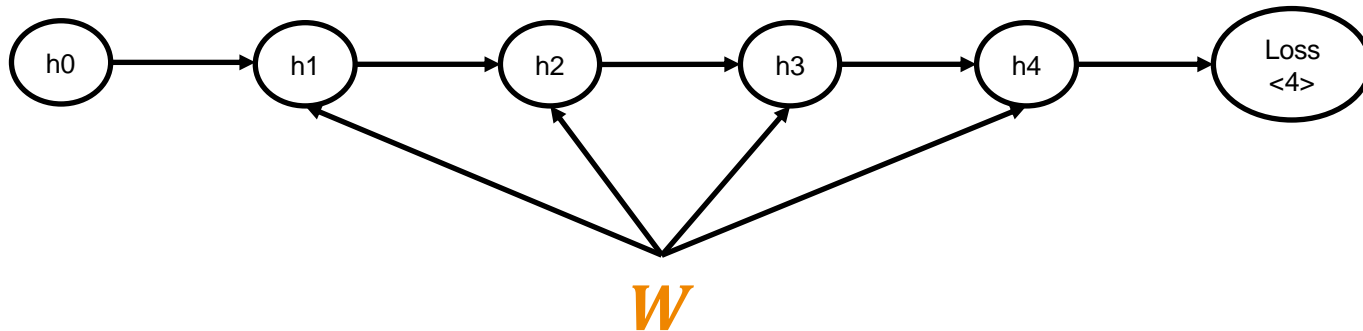
- Tìm mối liên kết giữa hàm $Loss$ và ma trận W

- $\partial Loss^{<4>}$ phụ thuộc vào h_4
 - h_4 phụ thuộc vào h_3 và W
 - h_3 phụ thuộc vào h_2 và W
 - h_2 phụ thuộc vào h_1 và W
 - h_1 phụ thuộc vào h_0 và W
- h_0 là hằng số - starting state



3 Backpropagation through time

Đạo hàm riêng hàm L trên biến W



- Áp dụng chain rule :

$$\begin{aligned}
 \frac{\partial Loss^{<4>}}{\partial W} &= \frac{\partial Loss^{<4>}}{\partial h_4} \frac{\partial h_4}{\partial W} \\
 &= [-y \log \hat{y} - (1 - y) \log(1 - \hat{y})]' \\
 &= [-y \log(\text{softmax}(Vh_t + b_y)) - (1 - y) \log(1 - \text{softmax}(Vh_t + b_y))]'
 \end{aligned}$$

3 Backpropagation through time

Đạo hàm riêng hàm L trên biến W

Tại sao $\frac{\partial h_4}{\partial W}$ hard ☹?

- Nhắc lại công thức tính h : $h_4 = \tanh(Wh_3 + Ux_4 + b)$
- Bỏ qua activation function, đạo hàm của h_4 trên biến W là

$$\frac{\partial h_4}{\partial W} = h_3 \quad ???$$

- Trong khi $h_3 = (Wh_2 + Ux_3 + b)$ (W xuất hiện trong h_3 , vậy h_3 không thể là hằng số)

3 Backpropagation through time

Đạo hàm riêng hàm L trên biến W

Tại sao $\frac{\partial h_4}{\partial W}$ hard ☹?

- Nhắc lại công thức tính h : $h_4 = (Wh_3 + Ux_4 + b)$
- Không thể tính $\frac{\partial h_4}{\partial W}$ bằng cách xem h_3 là hằng số được vì h_3 cũng phụ thuộc vào W
- Để tính được $\frac{\partial h_4}{\partial W}$, ta chia đạo hàm thành 2 phần:
 - Explicit
 - Implicit

3 Backpropagation through time

Đạo hàm riêng hàm L trên biến W

$$\begin{aligned}
 \frac{\partial h_4}{\partial W} &= \frac{\partial^+ h_4}{\partial W} + \frac{\partial h_4}{\partial h_3} \frac{\partial h_3}{\partial W} \\
 &= \frac{\partial^+ h_4}{\partial W} + \frac{\partial h_4}{\partial h_3} \left[\frac{\partial^+ h_3}{\partial W} + \frac{\partial h_3}{\partial h_2} \frac{\partial h_2}{\partial W} \right] \\
 &= \frac{\partial^+ h_4}{\partial W} + \frac{\partial h_4}{\partial h_3} \frac{\partial^+ h_3}{\partial W} + \frac{\partial h_4}{\partial h_3} \frac{\partial h_3}{\partial h_2} \left[\frac{\partial^+ h_2}{\partial W} + \frac{\partial h_2}{\partial h_1} \frac{\partial h_1}{\partial W} \right] \\
 &= \frac{\partial^+ h_4}{\partial W} + \frac{\partial h_4}{\partial h_3} \frac{\partial^+ h_3}{\partial W} + \frac{\partial h_4}{\partial h_3} \frac{\partial h_3}{\partial h_2} \frac{\partial^+ h_2}{\partial W} + \frac{\partial h_4}{\partial h_3} \frac{\partial h_3}{\partial h_2} \frac{\partial h_2}{\partial h_1} \left[\frac{\partial^+ h_1}{\partial W} \right] \\
 \frac{\partial h_4}{\partial W} &= \frac{\partial h_4}{\partial h_4} \frac{\partial^+ h_4}{\partial W} + \frac{\partial h_4}{\partial h_3} \frac{\partial^+ h_3}{\partial W} + \frac{\partial h_4}{\partial h_2} \frac{\partial^+ h_2}{\partial W} + \frac{\partial h_4}{\partial h_1} \left[\frac{\partial^+ h_1}{\partial W} \right] = \sum_{k=1}^4 \frac{\partial h_4}{\partial h_k} \frac{\partial^+ h_k}{\partial W}
 \end{aligned}$$

3 Backpropagation through time

Đạo hàm riêng hàm L trên biến W

$$\frac{\partial h_4}{\partial W} = \sum_{k=1}^4 \frac{\partial h_4}{\partial h_k} \frac{\partial^+ h_k}{\partial W}$$

- Quay lại bài toán cần tìm



Easy ☺

$$\frac{\partial Loss^{<4>}}{\partial W} = \frac{\partial Loss^{<4>}}{\partial h_4} \frac{\partial h_4}{\partial W}$$

Hard ☹

Tổng quát hóa:

$$\frac{\partial Loss^{<t>}}{\partial W} = \frac{\partial Loss^{<t>}}{\partial h_t} \sum_{k=1}^t \frac{\partial h_t}{\partial h_k} \frac{\partial^+ h_k}{\partial W}$$

$$\frac{\partial L}{\partial W} = \sum_{t=1}^T \frac{\partial Loss^{<t>}}{\partial W}$$

3

Backpropagation through time

Đạo hàm riêng hàm L trên biến U

Nội dung

1. Cấu trúc Recurrent Neural Network (RNN)
2. Cách hoạt động
3. Backpropagation through time
4. **Vanishing/Exploding gradients**

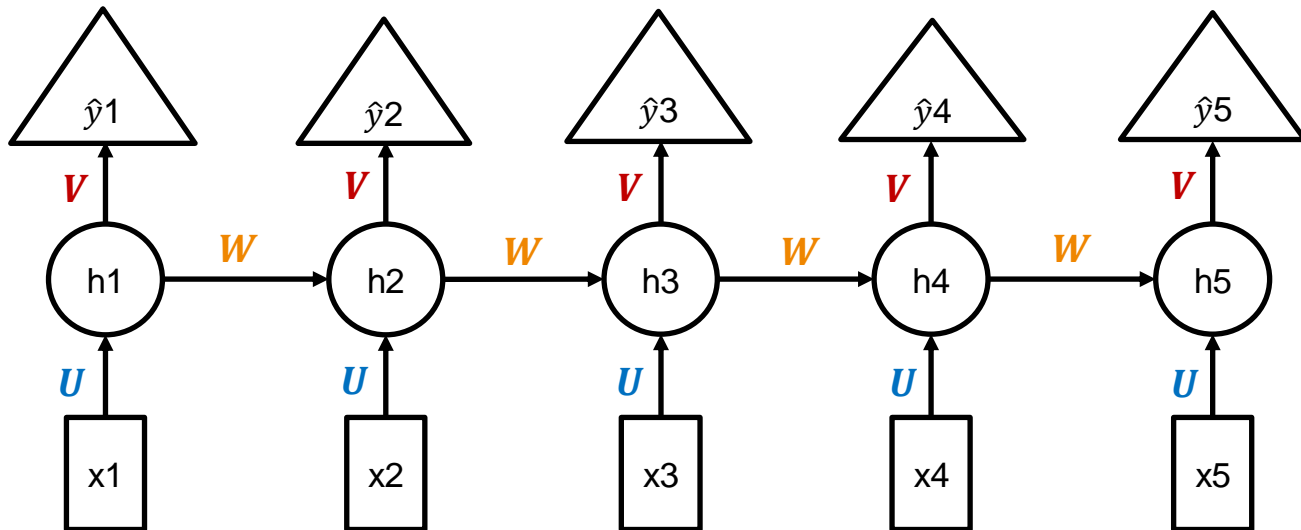
4 Long-term dependency

“In France, I had a great time and I learnt some of the _____ language.”



our parameters are not trained to capture long-term dependencies, so the word we predict will mostly depend on the previous few words, not much earlier ones

4 Vanishing/Exploding gradients



4 Vanishing/Exploding gradients

- We are interested in $\frac{\partial s_j}{\partial s_{j-1}}$

$$a_j = [a_{j1}, a_{j2}, a_{j3}, \dots, a_{jd}]$$

$$s_j = [\sigma(a_{j1}), \sigma(a_{j2}), \dots, \sigma(a_{jd})]$$

$$a_j = W s_{j-1} + b$$

$$s_j = \sigma(a_j)$$

$$\frac{\partial s_j}{\partial a_j} = \begin{bmatrix} \frac{\partial s_{j1}}{\partial a_{j1}} & \frac{\partial s_{j2}}{\partial a_{j1}} & \frac{\partial s_{j3}}{\partial a_{j1}} & \dots \\ \frac{\partial s_{j1}}{\partial a_{j2}} & \frac{\partial s_{j2}}{\partial a_{j2}} & \ddots & \\ \vdots & \vdots & \vdots & \frac{\partial s_{jd}}{\partial a_{jd}} \end{bmatrix}$$

$$= \begin{bmatrix} \sigma'(a_{j1}) & 0 & 0 & 0 \\ 0 & \sigma'(a_{j2}) & 0 & 0 \\ 0 & 0 & \ddots & \\ 0 & 0 & \dots & \sigma'(a_{jd}) \end{bmatrix}$$

$$= \text{diag}(\sigma'(a_j))$$

$$\frac{\partial s_j}{\partial s_{j-1}} = \frac{\partial s_j}{\partial a_j} \frac{\partial a_j}{\partial s_{j-1}}$$

$$= \text{diag}(\sigma'(a_j)) W$$

- We are interested in the magnitude of $\frac{\partial s_j}{\partial s_{j-1}} \leftarrow$ if it is small (large) $\frac{\partial s_t}{\partial s_k}$ and hence $\frac{\partial \mathcal{L}_t}{\partial W}$ will vanish (explode)

4 Vanishing/Exploding gradients

$$\begin{aligned}\left\| \frac{\partial s_j}{\partial s_{j-1}} \right\| &= \left\| \text{diag}(\sigma'(a_j))W \right\| \\ &\leq \left\| \text{diag}(\sigma'(a_j)) \right\| \|W\|\end{aligned}$$

$\because \sigma(a_j)$ is a bounded function (sigmoid, tanh) $\sigma'(a_j)$ is bounded

$$\begin{aligned}\sigma'(a_j) &\leq \frac{1}{4} = \gamma \text{ [if } \sigma \text{ is logistic]} \\ &\leq 1 = \gamma \text{ [if } \sigma \text{ is tanh]}\end{aligned}$$

$$\begin{aligned}\left\| \frac{\partial s_j}{\partial s_{j-1}} \right\| &\leq \gamma \|W\| \\ &\leq \gamma \lambda\end{aligned}$$

$$\begin{aligned}\left\| \frac{\partial s_t}{\partial s_k} \right\| &= \left\| \prod_{j=k+1}^t \frac{\partial s_j}{\partial s_{j-1}} \right\| \\ &\leq \prod_{j=k+1}^t \gamma \lambda \\ &\leq (\gamma \lambda)^{t-k}\end{aligned}$$

- If $\gamma \lambda < 1$ the gradient will vanish
- If $\gamma \lambda > 1$ the gradient could explode

4

Vanishing/Exploding gradients

- Dùng activation function khác
- Khởi tạo lại ma trận trọng số
- Dùng các biến thể của RNN
 - Long-short term memory
 - Gated Recurrent Units