

Backpropagation Through Time in Detail

Mục 9.7.2 – Phân tích chi tiết

Nhóm trình bày

Ngày 21 tháng 12 năm 2025

- Trình bày Backpropagation Through Time (BPTT) trên RNN tuyến tính, tập trung vào công thức gradient.
- Làm rõ cách tính $\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{qh}}$, $\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{hx}}$, $\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{hh}}$ cùng các bias liên quan.
- Nhấn mạnh ý nghĩa của từng gradient trong cập nhật tham số mô hình.

Mô hình RNN tuyến tính

- Chuỗi đầu vào $\mathbf{x}_t \in \mathbb{R}^d$, hidden state $\mathbf{h}_t \in \mathbb{R}^h$, đầu ra $\mathbf{o}_t \in \mathbb{R}^q$.
- Tham số: $\mathbf{W}_{hx} \in \mathbb{R}^{h \times d}$, $\mathbf{W}_{hh} \in \mathbb{R}^{h \times h}$, $\mathbf{W}_{qh} \in \mathbb{R}^{q \times h}$.
- Bias: $\mathbf{b}_h \in \mathbb{R}^h$ cho hidden layer, $\mathbf{b}_q \in \mathbb{R}^q$ cho output layer.
- computational graph mở dọc thời gian, tất cả bước t dùng chung các trọng số và bias.

Phương trình forward và loss

$$\mathbf{h}_t = \mathbf{W}_{hx}\mathbf{x}_t + \mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{b}_h, \quad (1)$$

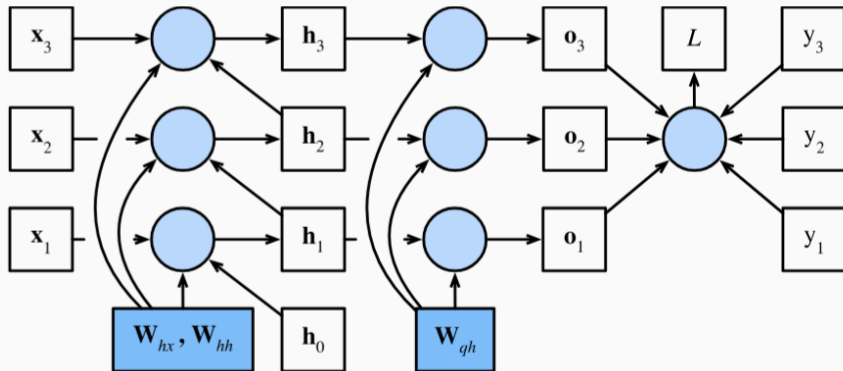
$$\mathbf{o}_t = \mathbf{W}_{qh}\mathbf{h}_t + \mathbf{b}_q \quad (2)$$

$$\mathcal{L} = \frac{1}{T} \sum_{t=1}^T \ell(\mathbf{o}_t, \mathbf{y}_t) \quad (3)$$

- Hidden state nhận thông tin từ input hiện tại và hidden state trước đó.
- Loss trung bình giúp gradient không phụ thuộc độ dài chuỗi.

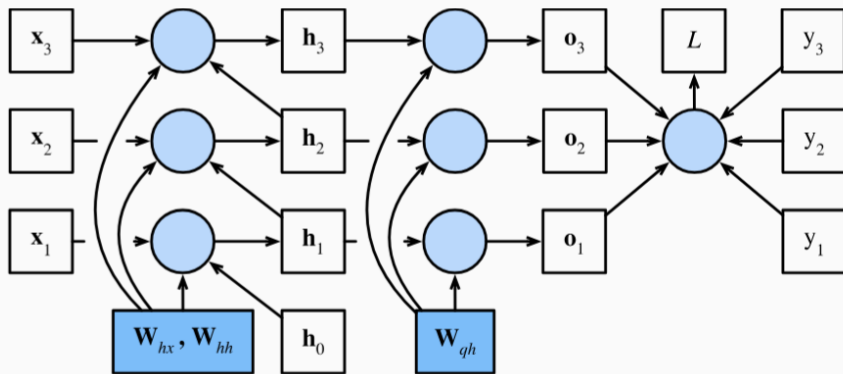
Đồ thị phụ thuộc

Để trực quan hóa phụ thuộc giữa biến và tham số trong quá trình tính toán RNN, ta vẽ một computational graph như Hình 9.7.2. Ví dụ, việc tính hidden state tại bước thời gian 3, \mathbf{h}_3 , phụ thuộc vào các tham số \mathbf{W}_{hx} và \mathbf{W}_{hh} , hidden state của bước trước \mathbf{h}_2 , và input tại bước hiện tại \mathbf{x}_3 .



Đồ thị phụ thuộc (tiếp)

Theo các phụ thuộc trong Hình 9.7.2, ta có thể đi ngược chiều mũi tên để lần lượt tính và lưu các gradient.



Gradient theo output tại time step t

$$L = \frac{1}{T} \sum_{k=1}^T \ell(o_k, y_k). \quad (1)$$

Vì o_t chỉ xuất hiện trong một hạng tử của tổng, ta có:

$$\begin{aligned} \frac{\partial L}{\partial o_t} &= \frac{\partial}{\partial o_t} \left(\frac{1}{T} \sum_{k=1}^T \ell(o_k, y_k) \right) \\ &= \frac{1}{T} \frac{\partial \ell(o_t, y_t)}{\partial o_t} \in \mathbb{R}^{q \times 1}. \end{aligned} \quad (2)$$

Gradient theo W_{qh} (1/2)

Output layer tuyến tính:

$$o_t = W_{qh}h_t + b_q. \quad (3)$$

Viết theo từng phần tử:

$$(o_t)_j = \sum_{i=1}^h (W_{qh})_{ji} (h_t)_i + (b_q)_j. \quad (4)$$

Áp dụng chain rule:

$$\begin{aligned} \frac{\partial L}{\partial (W_{qh})_{ji}} &= \sum_{t=1}^T \sum_{k=1}^q \frac{\partial L}{\partial (o_t)_k} \frac{\partial (o_t)_k}{\partial (W_{qh})_{ji}} \\ &= \sum_{t=1}^T \frac{\partial L}{\partial (o_t)_j} (h_t)_i. \end{aligned} \quad (5)$$

Gradient theo W_{qh} (2/2)

Viết dưới dạng ma trận:

$$\frac{\partial L}{\partial W_{qh}} = \sum_{t=1}^T \left(\frac{\partial L}{\partial o_t} \right) h_t^\top \in \mathbb{R}^{q \times h}. \quad (6)$$

- Gradient cộng dồn theo thời gian vì W_{qh} được chia sẻ cho mọi bước.

Gradient theo h_T

Tại thời điểm cuối:

$$o_T = W_{qh}h_T + b_q. \quad (7)$$

Áp dụng chain rule theo chỉ số:

$$\begin{aligned} \frac{\partial L}{\partial (h_T)_i} &= \sum_{j=1}^q \frac{\partial L}{\partial (o_T)_j} \frac{\partial (o_T)_j}{\partial (h_T)_i} \\ &= \sum_{j=1}^q (W_{qh})_{ji} \frac{\partial L}{\partial (o_T)_j}. \end{aligned} \quad (8)$$

Viết dạng vector:

$$\frac{\partial L}{\partial h_T} = W_{qh}^\top \frac{\partial L}{\partial o_T} \in \mathbb{R}^{h \times 1}. \quad (9)$$

Gradient theo h_t với $t < T$

Forward:

$$h_{t+1} = W_{hx}x_{t+1} + W_{hh}h_t + b_h, \quad o_t = W_{qh}h_t + b_q. \quad (10)$$

Áp dụng chain rule:

$$\frac{\partial L}{\partial h_t} = \left(\frac{\partial h_{t+1}}{\partial h_t} \right)^\top \frac{\partial L}{\partial h_{t+1}} + \left(\frac{\partial o_t}{\partial h_t} \right)^\top \frac{\partial L}{\partial o_t}. \quad (11)$$

Vì các ánh xạ là tuyến tính:

$$\frac{\partial h_{t+1}}{\partial h_t} = W_{hh}, \quad \frac{\partial o_t}{\partial h_t} = W_{qh}. \quad (12)$$

Suy ra:

$$\frac{\partial L}{\partial h_t} = W_{hh}^\top \frac{\partial L}{\partial h_{t+1}} + W_{qh}^\top \frac{\partial L}{\partial o_t}. \quad (13)$$

Gradient theo W_{hx} (1/2)

Hidden state tại thời điểm t :

$$h_t = W_{hx}x_t + W_{hh}h_{t-1} + b_h. \quad (14)$$

Vì L phụ thuộc vào W_{hx} thông qua toàn bộ các h_t , ta áp dụng chain rule:

$$\frac{\partial L}{\partial W_{hx}} = \sum_{t=1}^T \text{prod} \left(\frac{\partial L}{\partial h_t}, \frac{\partial h_t}{\partial W_{hx}} \right). \quad (15)$$

Xét từng phần tử của h_t :

$$(h_t)_i = \sum_{j=1}^d (W_{hx})_{ij}(x_t)_j + \sum_{k=1}^h (W_{hh})_{ik}(h_{t-1})_k + (b_h)_i. \quad (16)$$

Suy ra:

$$\frac{\partial (h_t)_i}{\partial (W_{hx})_{ij}} = (x_t)_j. \quad (17)$$

Gradient theo W_{hx} (2/2)

Do đó:

$$\begin{aligned}\frac{\partial L}{\partial (W_{hx})_{ij}} &= \sum_{t=1}^T \frac{\partial L}{\partial (h_t)_i} \frac{\partial (h_t)_i}{\partial (W_{hx})_{ij}} \\ &= \sum_{t=1}^T \frac{\partial L}{\partial (h_t)_i} (x_t)_j.\end{aligned}\tag{18}$$

Viết dưới dạng ma trận:

$$\frac{\partial L}{\partial W_{hx}} = \sum_{t=1}^T \left(\frac{\partial L}{\partial (h_t)_i} \right) x_t^\top \in \mathbb{R}^{h \times d}.\tag{19}$$

Gradient theo W_{hh} (1/2)

Từ phương trình hidden state:

$$h_t = W_{hx}x_t + W_{hh}h_{t-1} + b_h, \quad (20)$$

W_{hh} ảnh hưởng đến L thông qua toàn bộ chuỗi hidden states h_1, \dots, h_T .

Áp dụng chain rule:

$$\frac{\partial L}{\partial W_{hh}} = \sum_{t=1}^T \text{prod} \left(\frac{\partial L}{\partial h_t}, \frac{\partial h_t}{\partial W_{hh}} \right). \quad (21)$$

Viết theo từng phần tử:

$$(h_t)_i = \sum_{k=1}^h (W_{hh})_{ik} (h_{t-1})_k + \dots \quad (22)$$

Suy ra:

$$\frac{\partial (h_t)_i}{\partial (W_{hh})_{ik}} = (h_{t-1})_k. \quad (23)$$

Gradient theo W_{hh} (2/2)

Do đó:

$$\begin{aligned}\frac{\partial L}{\partial (W_{hh})_{ik}} &= \sum_{t=1}^T \frac{\partial L}{\partial (h_t)_i} \frac{\partial (h_t)_i}{\partial (W_{hh})_{ik}} \\ &= \sum_{t=1}^T \frac{\partial L}{\partial (h_t)_i} (h_{t-1})_k.\end{aligned}\tag{24}$$

Viết dưới dạng ma trận:

$$\frac{\partial L}{\partial W_{hh}} = \sum_{t=1}^T \left(\frac{\partial L}{\partial h_t} \right) h_{t-1}^\top \in \mathbb{R}^{h \times h}.\tag{25}$$

Gradient theo bias b_h

Hidden state:

$$h_t = W_{hx}x_t + W_{hh}h_{t-1} + b_h. \quad (26)$$

Xét đạo hàm theo b_h :

$$\frac{\partial h_t}{\partial b_h} = I_h. \quad (27)$$

Áp dụng chain rule:

$$\begin{aligned} \frac{\partial L}{\partial b_h} &= \sum_{t=1}^T \frac{\partial L}{\partial h_t} \frac{\partial h_t}{\partial b_h} \\ &= \sum_{t=1}^T \frac{\partial L}{\partial h_t} \in \mathbb{R}^{h \times 1}. \end{aligned} \quad (28)$$

Gradient theo b_q

Từ biểu thức:

$$o_t = W_{qh}h_t + b_q, \quad (29)$$

ta có:

$$\frac{\partial o_t}{\partial b_q} = I_q. \quad (30)$$

Áp dụng chain rule:

$$\begin{aligned} \frac{\partial L}{\partial b_q} &= \sum_{t=1}^T \frac{\partial L}{\partial o_t} \frac{\partial o_t}{\partial b_q} \\ &= \sum_{t=1}^T \frac{\partial L}{\partial o_t} \in \mathbb{R}^{q \times 1}. \end{aligned} \quad (31)$$

Kết luận: Gradient theo các tham số của mô hình

Gradient theo output layer:

$$\begin{aligned}\frac{\partial L}{\partial W_{qh}} &= \sum_{t=1}^T \left(\frac{\partial L}{\partial o_t} \right) h_t^\top, \\ \frac{\partial L}{\partial b_q} &= \sum_{t=1}^T \frac{\partial L}{\partial o_t}.\end{aligned}\tag{32}$$

Gradient theo hidden layer:

$$\begin{aligned}\frac{\partial L}{\partial W_{hx}} &= \sum_{t=1}^T \left(\frac{\partial L}{\partial h_t} \right) x_t^\top, \\ \frac{\partial L}{\partial W_{hh}} &= \sum_{t=1}^T \left(\frac{\partial L}{\partial h_t} \right) h_{t-1}^\top, \\ \frac{\partial L}{\partial b_h} &= \sum_{t=1}^T \frac{\partial L}{\partial h_t}.\end{aligned}\tag{33}$$