

Recurrent Neural Networks and the Exploding and Vanishing Gradient Problems

Le Quang Trung, Nguyen Ngoc Hung, Tran Le Anh Pha

Ngày 5 tháng 1 năm 2026

Mục lục chương 1

1 Introduction & Notation

2 Train Recurrent Neural Networks

3 Exploding And Vanishing Gradient Problems

Sequential Data

Dữ liệu chuỗi (sequential data) là các phần tử được sắp theo thứ tự tạo thành chuỗi.

Ví dụ: chuỗi thời gian, chuỗi DNA, hoặc chuỗi hành vi người dùng.

Mô hình chuỗi (sequence models) nhận/ xuất dữ liệu dạng chuỗi.

What Are Recurrent Neural Networks?

Recurrent Neural Networks (RNNs) là kiến trúc mạng nơ-ron dùng để phát hiện quy luật trong dữ liệu chuỗi.

Có thể áp dụng cho ảnh sau khi tách ảnh thành các patch có thứ tự. Điểm khác biệt với feedforward neural networks (MLP) là cách thông tin được truyền qua thời gian.

RNN có chu trình (cycles) truyền thông tin ngược trong mạng, cho phép mô hình xét các input trước đó $X_{0:t-1}$ thay vì chỉ input hiện tại X_t .

Model Architecture

Gọi hidden state tại thời điểm t là H_t và input là X_t . Các tham số được chia sẻ: W_{xh} , W_{hh} , W_{ho} , cùng bias b_h , b_o . Activation function ϕ thường là logistic sigmoid hoặc tanh.

$$\begin{aligned}H_t &= \phi_h(X_t W_{xh} + H_{t-1} W_{hh} + b_h), \\O_t &= \phi_o(H_t W_{ho} + b_o).\end{aligned}$$

Loss Function

Hàm loss đánh giá hiệu năng bằng cách so sánh output y_t với target z_t :

$$\mathcal{L}(y, z) = \sum_{t=1}^T \ell_t(y_t, z_t).$$

Loss phụ thuộc bài toán; ví dụ: Euclidean distance, Hamming distance, cross-entropy.

Ưu điểm

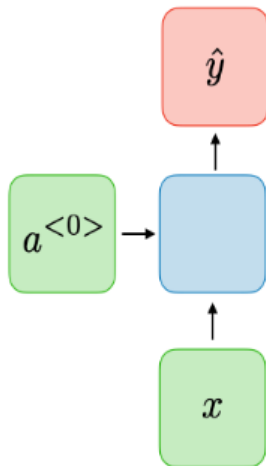
- Xử lý chuỗi độ dài bất kỳ, chia sẻ weights qua thời gian.
- Kích thước mô hình không tăng theo độ dài chuỗi.
- Tính toán xét đến thông tin quá khứ.

Nhược điểm

- Tính toán có thể chậm.
- Khó lấy thông tin ở quá khứ xa.
- Không xét được input tương lai khi tạo state hiện tại.

Various Usage of RNNs: One-to-One

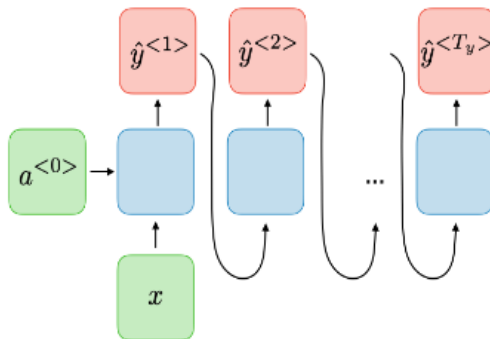
One-to-One RNNs (vanilla neural networks) dùng cho tác vụ một input, một output.



Various Usage of RNNs: One-to-Many

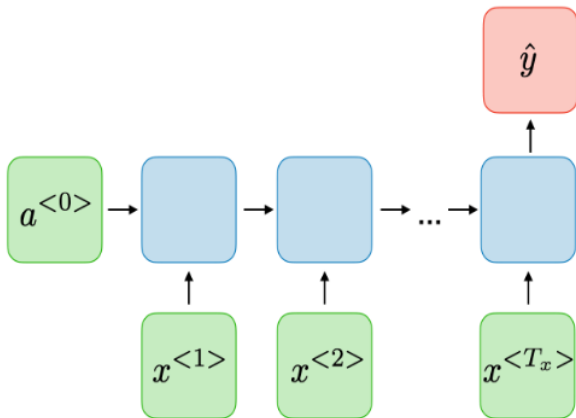
One-to-Many RNNs sinh nhiều output từ một input.

Ứng dụng: music generation, image captioning.



Various Usage of RNNs: Many-to-One

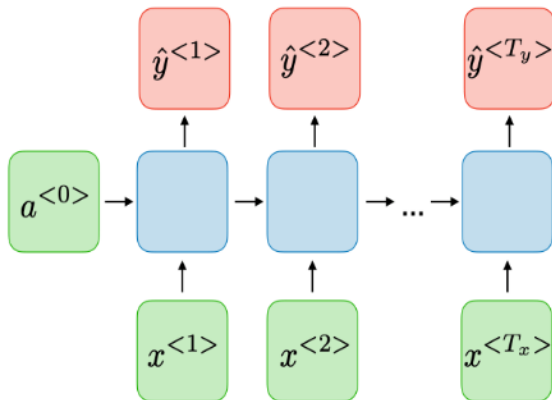
Many-to-One RNNs nhận chuỗi input và xuất một output cố định.
Ví dụ phổ biến: sentiment analysis.



Various Usage of RNNs: Many-to-Many

Many-to-Many RNNs ánh xạ chuỗi input sang chuỗi output.
Số bước bằng nhau: số input = số output (ví dụ: named-entity recognition).

Số bước khác nhau: độ dài input khác output (ví dụ: machine translation).



Mục lục chương 2

1 Introduction & Notation

2 Train Recurrent Neural Networks

3 Exploding And Vanishing Gradient Problems

Yêu cầu bổ sung: Các hình vẽ và kết quả demo cần được nhúng vào trong slides.

Thiết lập ký hiệu & công thức sẽ dùng

Input/hidden/output: \mathbf{x}_t (đầu vào), \mathbf{h}_t (hidden state), \mathbf{o}_t (đầu ra).

Tiền kích hoạt hidden: $\mathbf{a}_t = \mathbf{W}_{hx}\mathbf{x}_t + \mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{b}_h$.

Hidden: $\mathbf{h}_t = \phi_h(\mathbf{a}_t)$, đạo hàm ϕ'_h (ví dụ tanh).

Tiền kích hoạt output: $\mathbf{z}_t = \mathbf{W}_{qh}\mathbf{h}_t + \mathbf{b}_q$.

Output: $\mathbf{o}_t = \phi_o(\mathbf{z}_t)$, với ϕ_o là activation/output layer (tuyến tính, softmax, sigmoid... tùy bài toán).

Loss chuỗi: $\mathcal{L} = \frac{1}{T} \sum_{t=1}^T \ell(\mathbf{o}_t, \mathbf{y}_t)$.

Quy trình huấn luyện RNN

Quy trình thường sẽ như thế này:

Forward Pass: tính toán đầu ra và sai số (Loss).

Backward Pass: tính gradient ∇J .

Update Rule (Optimizer): dùng gradient để cập nhật tham số với các lựa chọn phổ biến:

- Batch Gradient Descent

- Stochastic Gradient Descent

- Mini-batch Gradient Descent

Forward pass: tính toán qua thời gian

Khởi tạo $\mathbf{h}_0 = \mathbf{0}$.

Với mỗi bước $t = 1, \dots, T$:

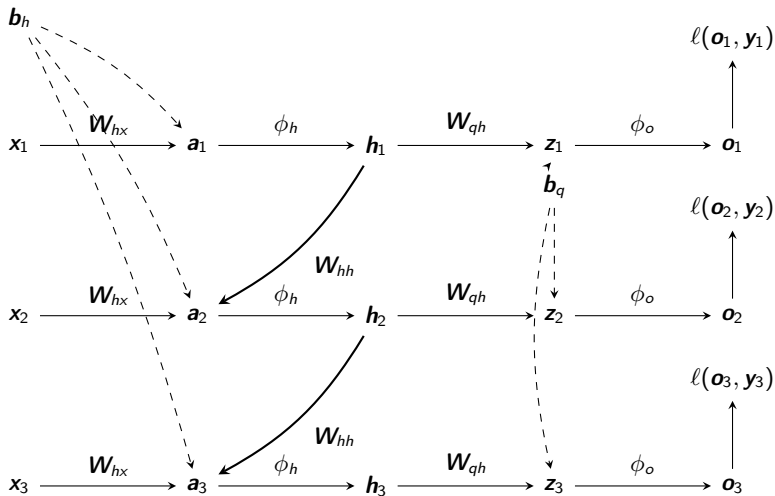
$$\mathbf{a}_t = \mathbf{W}_{hx}\mathbf{x}_t + \mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{b}_h, \quad \mathbf{h}_t = \phi_h(\mathbf{a}_t)$$

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

Tính loss tổng/ trung bình:

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

Computational graph



Hai insight quan trọng khi lấy đạo hàm theo tham số:

Weight Sharing: cùng một tham số W_{hx} , W_{hh} , W_{qh} xuất hiện ở mọi timestep, nên

$$\frac{\partial \mathcal{L}}{\partial W} = \sum_{t=1}^T \frac{\partial \mathcal{L}_t}{\partial W}.$$

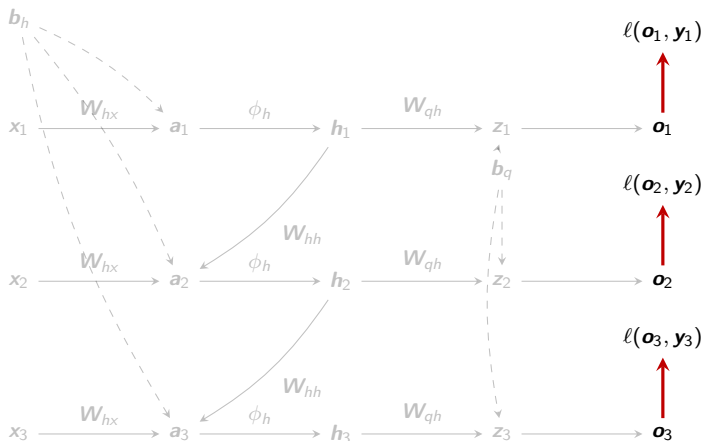
Phụ thuộc qua thời gian: gradient theo hidden tại thời điểm t nhận đóng góp từ loss hiện tại và từ tương lai:

$$\frac{\partial \mathcal{L}}{\partial h_t} = \frac{\partial \mathcal{L}_t}{\partial h_t} + W_{hh}^\top \frac{\partial \mathcal{L}}{\partial a_{t+1}}.$$

Đường đạo hàm cho \mathbf{o}_t

Mỗi \mathbf{o}_t chỉ nối với loss tại cùng thời điểm: $\mathbf{o}_t \rightarrow \ell_t$.

Không có cộng dồn qua các bước thời gian (chỉ nhân hệ số $\frac{1}{T}$).



Gradient theo output đầu ra \mathbf{o}_t

Loss chuỗi: $\mathcal{L} = \frac{1}{T} \sum_{k=1}^T \ell(\mathbf{o}_k, \mathbf{y}_k).$

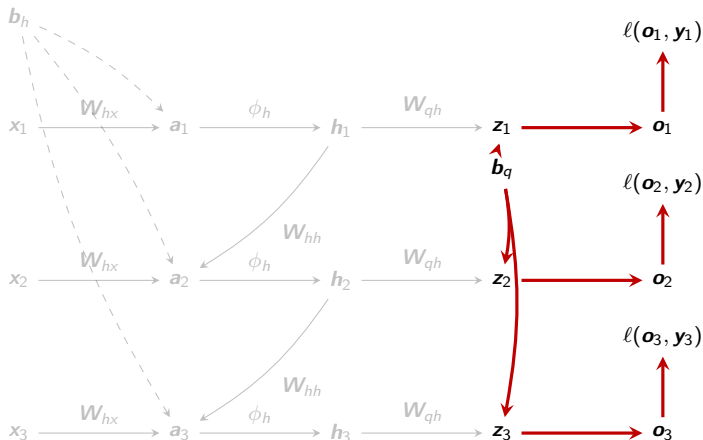
Mỗi thời điểm xuất hiện đúng một lần trong tổng nên:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{o}_t} = \frac{1}{T} \frac{\partial \ell(\mathbf{o}_t, \mathbf{y}_t)}{\partial \mathbf{o}_t} \in \mathbb{R}^q.$$

Đường đạo hàm cho \mathbf{b}_q

Gradient theo \mathbf{b}_q đi dọc nhánh $\mathbf{b}_q \rightarrow \mathbf{z}_t \rightarrow \mathbf{o}_t \rightarrow \ell_t$ tại từng bước thời gian.

Vì \mathbf{b}_q được chia sẻ, các đóng góp này được cộng lại rồi chia T .



Gradient theo bias đầu ra \mathbf{b}_q

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

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

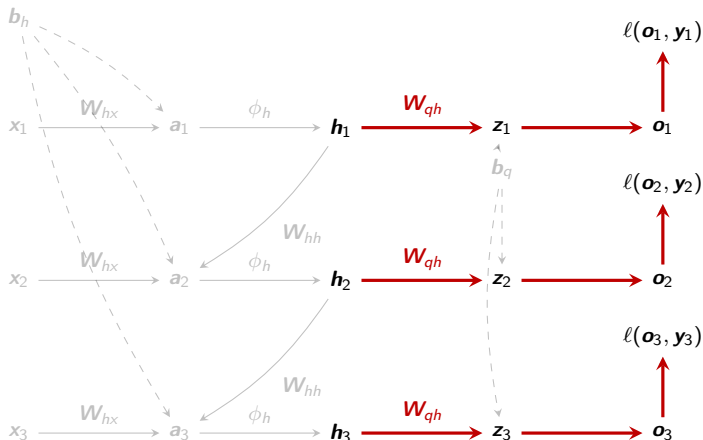
Vì \mathbf{b}_q dùng chung cho mọi bước thời gian:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \mathbf{b}_q} &= \frac{1}{T} \sum_{t=1}^T \frac{\partial \ell(\mathbf{o}_t, \mathbf{y}_t)}{\partial \mathbf{o}_t} \cdot \frac{\partial \mathbf{o}_t}{\partial \mathbf{z}_t} \cdot \frac{\partial \mathbf{z}_t}{\partial \mathbf{b}_q} \quad (\text{chuỗi phụ thuộc } \mathbf{b}_q \rightarrow \mathbf{z}_t \rightarrow \mathbf{o}_t \rightarrow \ell) \\ &= \frac{1}{T} \sum_{t=1}^T \frac{\partial \ell(\mathbf{o}_t, \mathbf{y}_t)}{\partial \mathbf{o}_t} \odot \phi'_o(\mathbf{z}_t) \quad (\text{do } \partial \mathbf{z}_t / \partial \mathbf{b}_q = \mathbf{I}) \end{aligned}$$

Đường đạo hàm cho \mathbf{W}_{qh}

Gradient đi theo nhánh $\mathbf{W}_{qh} \rightarrow \mathbf{z}_t \rightarrow \mathbf{o}_t \rightarrow \ell_t$, được nhân thêm \mathbf{h}_t^\top tại mỗi bước.

Vì \mathbf{W}_{qh} dùng chung, cộng các đóng góp qua thời gian rồi chia T .



Gradient theo trọng số \mathbf{W}_{qh}

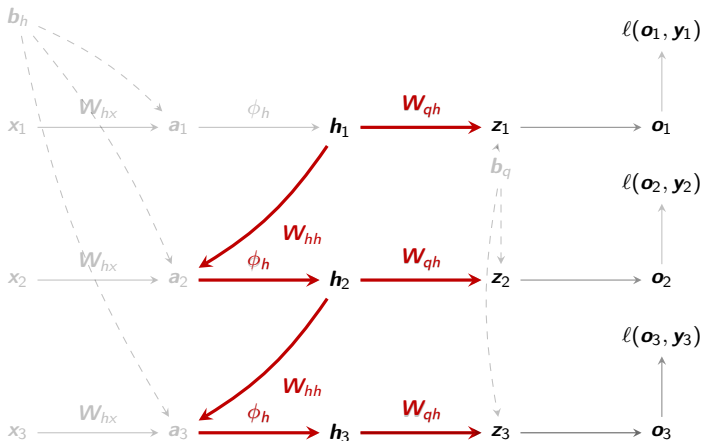
Sử dụng $\mathbf{z}_t = \mathbf{W}_{qh}\mathbf{h}_t + \mathbf{b}_q$ và áp dụng chain rule từng bước:

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{qh}} &= \frac{1}{T} \sum_{t=1}^T \frac{\partial \ell(\mathbf{o}_t, \mathbf{y}_t)}{\partial \mathbf{W}_{qh}} && \text{(tách tổng theo từng thời điểm)} \\ &= \frac{1}{T} \sum_{t=1}^T \frac{\partial \ell(\mathbf{o}_t, \mathbf{y}_t)}{\partial \mathbf{o}_t} \cdot \frac{\partial \mathbf{o}_t}{\partial \mathbf{z}_t} \cdot \frac{\partial \mathbf{z}_t}{\partial \mathbf{W}_{qh}} && \text{(chuỗi phụ thuộc } \mathbf{W}_{qh} \rightarrow \mathbf{z}_t \rightarrow \mathbf{o}_t \rightarrow \ell) \\ &= \frac{1}{T} \sum_{t=1}^T \left(\frac{\partial \ell(\mathbf{o}_t, \mathbf{y}_t)}{\partial \mathbf{o}_t} \odot \phi'_o(\mathbf{z}_t) \right) \mathbf{h}_t^\top\end{aligned}$$

Đường đạo hàm cho h_t

Hai nhánh: (1) qua output hiện tại $h_t \rightarrow z_t \rightarrow o_t \rightarrow \ell_t$; (2) hồi quy về tương lai $h_t \rightarrow a_{t+1} \rightarrow h_{t+1} \rightarrow \dots$.

Tổng hợp cả hai nhánh (nhánh tương lai gói trong $\frac{\partial \mathcal{L}}{\partial h_{t+1}}$).



Gradient theo hidden state \mathbf{h}_t

$$\frac{\partial \mathcal{L}}{\partial \mathbf{h}_t} = \underbrace{\frac{1}{T} \frac{\partial \ell(\mathbf{o}_t, \mathbf{y}_t)}{\partial \mathbf{h}_t}}_{\text{nhánh output tại } t} + \underbrace{\frac{\partial \mathcal{L}}{\partial \mathbf{h}_{t+1}} \cdot \frac{\partial \mathbf{h}_{t+1}}{\partial \mathbf{h}_t}}_{\text{nhánh hồi quy về tương lai}}.$$

Mở từng nhánh bằng chain rule:

$$\frac{\partial \ell(\mathbf{o}_t, \mathbf{y}_t)}{\partial \mathbf{h}_t} = \frac{\partial \ell(\mathbf{o}_t, \mathbf{y}_t)}{\partial \mathbf{o}_t} \cdot \frac{\partial \mathbf{o}_t}{\partial \mathbf{z}_t} \cdot \frac{\partial \mathbf{z}_t}{\partial \mathbf{h}_t} \quad (\text{chuỗi } \mathbf{h}_t \rightarrow \mathbf{z}_t \rightarrow \mathbf{o}_t \rightarrow \ell)$$

$$= \mathbf{W}_{qh}^\top \left(\frac{\partial \ell(\mathbf{o}_t, \mathbf{y}_t)}{\partial \mathbf{o}_t} \odot \phi'_o(\mathbf{z}_t) \right),$$

$$\frac{\partial \mathbf{h}_{t+1}}{\partial \mathbf{h}_t} = \frac{\partial \phi_h(\mathbf{a}_{t+1})}{\partial \mathbf{a}_{t+1}} \cdot \frac{\partial \mathbf{a}_{t+1}}{\partial \mathbf{h}_t} = \text{diag}(\phi'_h(\mathbf{a}_{t+1})) \mathbf{W}_{hh}^\top.$$

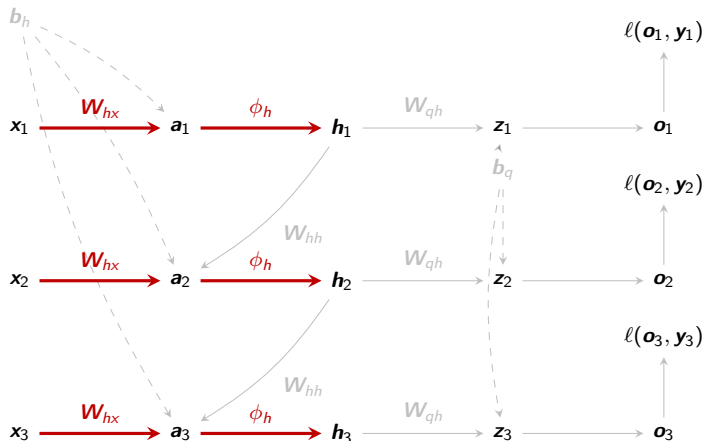
Ghép lại:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \mathbf{h}_t} &= \frac{1}{T} \mathbf{W}_{qh}^\top \left(\frac{\partial \ell}{\partial \mathbf{o}_t} \odot \phi'_o(\mathbf{z}_t) \right) \\ &\quad + \mathbf{W}_{hh}^\top \left(\frac{\partial \mathcal{L}}{\partial \mathbf{h}_{t+1}} \odot \phi'_h(\mathbf{a}_{t+1}) \right), \end{aligned}$$

Đường đạo hàm cho W_{hx}

Chuỗi: $W_{hx} \rightarrow a_t \rightarrow h_t \rightarrow z_t \rightarrow o_t \rightarrow \ell_t$ tại mỗi timestep.

Cộng các đóng góp qua thời gian (do chia sẻ tham số) rồi lấy trung bình.



Gradient theo \mathbf{W}_{hx}

Nhắc lại tiền kích hoạt & hidden:

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

Với \mathbf{W}_{hx} (input \rightarrow hidden):

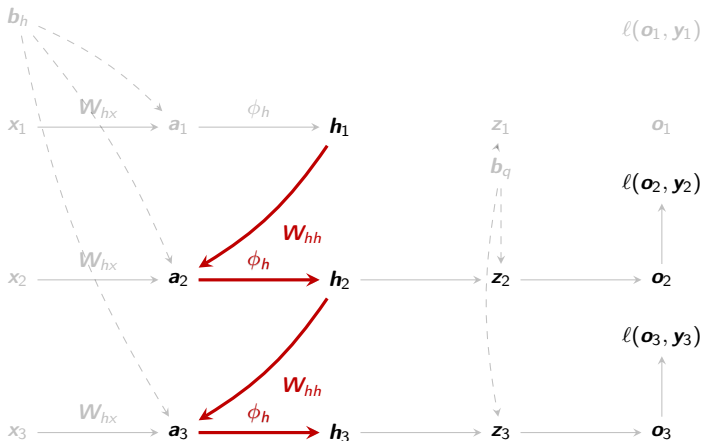
$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \mathbf{W}_{hx}} &= \sum_{t=1}^T \frac{\partial \mathcal{L}}{\partial \mathbf{h}_t} \cdot \frac{\partial \mathbf{h}_t}{\partial \mathbf{a}_t} \cdot \frac{\partial \mathbf{a}_t}{\partial \mathbf{W}_{hx}} && (\text{dùng trực tiếp } \frac{\partial \mathcal{L}}{\partial \mathbf{h}_t} \text{ đã suy ra}) \\ &= \sum_{t=1}^T \left(\frac{\partial \mathcal{L}}{\partial \mathbf{h}_t} \odot \phi'_h(\mathbf{a}_t) \right) \mathbf{x}_t^\top, \end{aligned}$$

vì $\partial \mathbf{a}_t / \partial \mathbf{W}_{hx} = \mathbf{x}_t$. Lưu ý $\frac{\partial \mathcal{L}}{\partial \mathbf{h}_t}$ đã bao gồm hệ số $1/T$ và các đóng góp hồi quy phía sau.

Đường đạo hàm cho W_{hh}

Chuỗi: $\tilde{h}_{t-1} \xrightarrow{W_{hh}} a_t \rightarrow h_t \rightarrow z_t \rightarrow o_t \rightarrow \ell_t$.

Cộng dồn qua các timestep có kết nối hồi quy, rồi chia trung bình.



Gradient theo \mathbf{W}_{hh}

Với $\mathbf{a}_t = \mathbf{W}_{hx}\mathbf{x}_t + \mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{b}_h$ suy ra $\partial\mathbf{a}_t/\partial\mathbf{W}_{hh} = \mathbf{h}_{t-1}$.

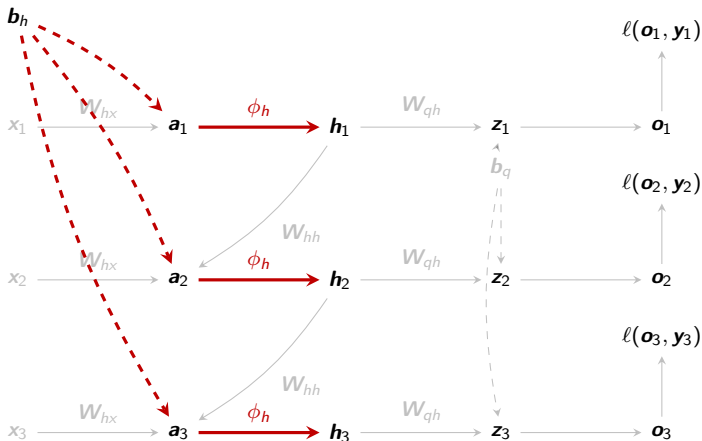
$$\begin{aligned}\frac{\partial\mathcal{L}}{\partial\mathbf{W}_{hh}} &= \sum_{t=1}^T \frac{\partial\mathcal{L}}{\partial\mathbf{h}_t} \cdot \frac{\partial\mathbf{h}_t}{\partial\mathbf{a}_t} \cdot \frac{\partial\mathbf{a}_t}{\partial\mathbf{W}_{hh}} && \text{(chain rule, dùng } \frac{\partial\mathcal{L}}{\partial\mathbf{h}_t} \text{ từ slide trước)} \\ &= \sum_{t=1}^T \left(\frac{\partial\mathcal{L}}{\partial\mathbf{h}_t} \odot \phi'_h(\mathbf{a}_t) \right) \mathbf{h}_{t-1}^\top,\end{aligned}$$

vì $\partial\mathbf{a}_t/\partial\mathbf{W}_{hh} = \mathbf{h}_{t-1}$. Hệ số $1/T$ và ảnh hưởng từ các bước tương lai đã nằm trong $\frac{\partial\mathcal{L}}{\partial\mathbf{h}_t}$.

Đường đạo hàm cho b_h

Nhánh: $b_h \dashrightarrow a_t \rightarrow h_t \rightarrow z_t \rightarrow o_t \rightarrow \ell_t$ tại mỗi t .

Bias chia sẻ nên các đóng góp được cộng lại rồi chia T .



Gradient theo bias \mathbf{b}_h

Từ $\mathbf{a}_t = \mathbf{W}_{hx}\mathbf{x}_t + \mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{b}_h$ ta có $\partial\mathbf{a}_t/\partial\mathbf{b}_h = \mathbf{I}$.

$$\frac{\partial\mathcal{L}}{\partial\mathbf{b}_h} = \sum_{t=1}^T \frac{\partial\mathcal{L}}{\partial\mathbf{h}_t} \cdot \frac{\partial\mathbf{h}_t}{\partial\mathbf{a}_t} \cdot \frac{\partial\mathbf{a}_t}{\partial\mathbf{b}_h} \quad (\text{chain rule qua } \mathbf{a}_t)$$

$$= \sum_{t=1}^T \left(\frac{\partial\mathcal{L}}{\partial\mathbf{h}_t} \odot \phi'_h(\mathbf{a}_t) \right) \quad (\text{do } \partial\mathbf{a}_t/\partial\mathbf{b}_h = \mathbf{I}).$$

Tương tự như \mathbf{W}_{hx} , hệ số $1/T$ và các ảnh hưởng lan ngược đã nằm trong $\frac{\partial\mathcal{L}}{\partial\mathbf{h}_t}$.

Tóm tắt gradient các tham số

$$\frac{\partial \mathcal{L}}{\partial \mathbf{b}_q} = \frac{1}{T} \sum_{t=1}^T \frac{\partial \ell(\mathbf{o}_t, \mathbf{y}_t)}{\partial \mathbf{o}_t} \odot \phi'_o(\mathbf{z}_t),$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{qh}} = \frac{1}{T} \sum_{t=1}^T \left(\frac{\partial \ell(\mathbf{o}_t, \mathbf{y}_t)}{\partial \mathbf{o}_t} \odot \phi'_o(\mathbf{z}_t) \right) \mathbf{h}_t^\top,$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{hx}} = \sum_{t=1}^T \left(\frac{\partial \mathcal{L}}{\partial \mathbf{h}_t} \odot \phi'_h(\mathbf{a}_t) \right) \mathbf{x}_t^\top,$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{hh}} = \sum_{t=1}^T \left(\frac{\partial \mathcal{L}}{\partial \mathbf{h}_t} \odot \phi'_h(\mathbf{a}_t) \right) \mathbf{h}_{t-1}^\top,$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{b}_h} = \sum_{t=1}^T \left(\frac{\partial \mathcal{L}}{\partial \mathbf{h}_t} \odot \phi'_h(\mathbf{a}_t) \right).$$

Điều kiện biên: $\frac{\partial \mathcal{L}}{\partial \mathbf{h}_{T+1}} = 0$.

Demo: RNN dự đoán giá cổ phiếu

Bài toán & kiến trúc

Many-to-one: cửa sổ 20 giá \rightarrow dự báo 1 giá tiếp theo.

Kiến trúc: 1 input, 16 hidden (tanh), 1 output tuyến tính.

Loss: MSE trên output cuối chuỗi; dữ liệu min-max fit trên train, dùng lại cho test.

Huấn luyện

300 epochs, $lr = 0.01$.

Gradient clipping $\tau = 5$ để tránh exploding.

Hiệu năng test

RMSE ≈ 1.45

MAE ≈ 1.19

Ghi chú

Chuỗi synthetic: sóng sin + nhiễu, 320 điểm train / 100 điểm test.

Dự báo dùng tham số đã huấn luyện; không dùng thêm thông tin test.

Code forward RNN

```
1 def forward(self, x_seq, h0=None, return_cache: bool = False):
2     # Forward cho mot chuoi (T, input_size). Tra outputs (T,
3     #   output_size) va h_T.
4     x_seq = np.asarray(x_seq, dtype=float)
5     if x_seq.ndim == 1:
6         x_seq = x_seq.reshape(-1, self.input_size)
7     T = x_seq.shape[0]
8     h_prev = np.zeros(self.hidden_size) if h0 is None else h0
9
10    hs = [h_prev]
11    pre_acts = []
12    outputs = []
13
14    for t in range(T):
15        a_t = x_seq[t].dot(self.W_hx) + h_prev.dot(self.W_hh)
16            + self.b_h
17        h_t = np.tanh(a_t)
18        o_t = h_t.dot(self.W_qh) + self.b_q
```

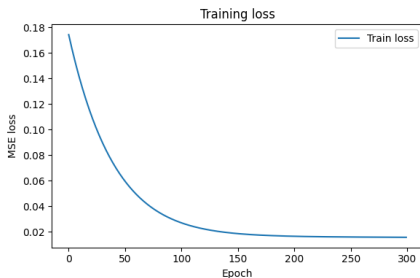
Code backward (1/2)

```
1 def backward(self, cache: dict, target: np.ndarray):
2     # BPTT cho loss  $0.5*(o_T - y)^2$ . Chỉ dùng output cuối
      chuoì.
3     x_seq = cache["x_seq"]
4     hs = cache["hs"]
5     pre_acts = cache["pre_acts"]
6     outputs = cache["outputs"]
7
8     T = x_seq.shape[0]
9     target = np.asarray(target).reshape(-1)
10    o_T = outputs[-1].reshape(-1)
11
12    grads = {
13        "W_hx": np.zeros_like(self.W_hx),
14        "W_hh": np.zeros_like(self.W_hh),
15        "W_qh": np.zeros_like(self.W_qh),
16        "b_h": np.zeros_like(self.b_h),
17        "b_q": np.zeros_like(self.b_q),
18    }
19
20    dL_do = o_T - target # derivative of  $0.5*(o - y)^2$ 
21    grads["W_qh"] += np.outer(hs[-1], dL_do)
```

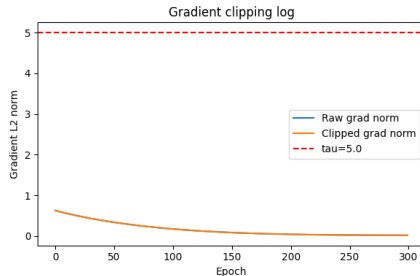
Code backward (2/2)

```
1     grads["b_q"] += dL_do
2
3     d_next = dL_do.dot(self.W_qh.T)
4
5     for t in reversed(range(T)):
6         h_t = hs[t + 1]
7         h_prev = hs[t]
8         a_t = pre_acts[t]
9         x_t = x_seq[t]
10
11         dh = d_next
12         da = dh * (1.0 - h_t**2)  # tanh'
13
14         grads["W_hx"] += np.outer(x_t, da)
15         grads["W_hh"] += np.outer(h_prev, da)
16         grads["b_h"] += da
17
18         d_next = da.dot(self.W_hh.T)
19
20     loss = 0.5 * np.mean((o_T - target) ** 2)
21     return loss, grads
```

Training loss và norm gradient (đã clipping)



Loss giảm đều trong 300 epochs.

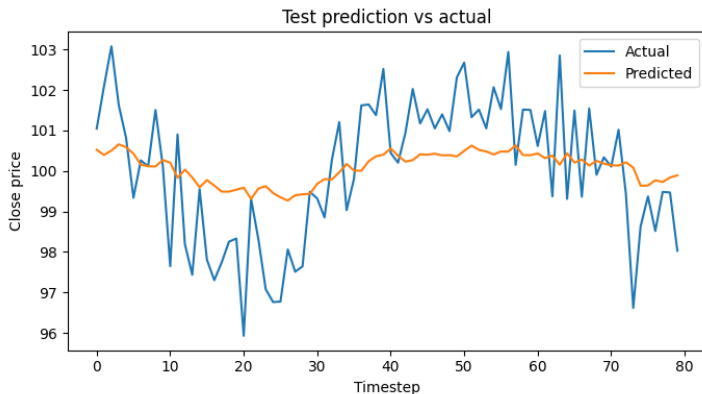


Clipping giữ gradient norm gần ngưỡng $\tau = 5$ và tránh exploding.

Dự báo trên tập test

Dùng tham số đã huấn luyện và scaler của train set để dự báo giá test.

Sai số: RMSE ≈ 1.45 , MAE ≈ 1.19 .



Mục lục chương 3

1 Introduction & Notation

2 Train Recurrent Neural Networks

3 Exploding And Vanishing Gradient Problems

Nguồn gốc của vanishing và exploding gradients

Từ chương trước, gradient theo hidden state tại thời điểm t có dạng:

$$\frac{\partial L}{\partial h_t} = \sum_{k=t}^T (W_{hh}^\top)^{k-t} W_{qh}^\top \frac{\partial L}{\partial o_k}. \quad (1)$$

Do đó, độ lớn của gradient phụ thuộc trực tiếp vào các lũy thừa của ma trận W_{hh}^\top .

Vấn đề cốt lõi:

Gradient bị nhân lặp lại bởi cùng một ma trận qua nhiều bước thời gian.

Hành vi của gradient được quyết định bởi W_{hh} (độ lớn/đặc trưng riêng của ma trận này).

Ước lượng độ lớn gradient

Lấy chuẩn hai vế:

$$\left\| \frac{\partial L}{\partial h_t} \right\| \leq \sum_{k=t}^T \left\| (W_{hh}^\top)^{k-t} \right\| \left\| W_{qh}^\top \right\| \left\| \frac{\partial L}{\partial o_k} \right\|. \quad (2)$$

Với chuẩn ma trận dưới chuẩn vector:

$$\left\| (W_{hh}^\top)^{k-t} \right\| \leq \|W_{hh}\|^{k-t}. \quad (3)$$

Do đó:

$$\left\| \frac{\partial L}{\partial h_t} \right\| \lesssim \sum_{k=t}^T \|W_{hh}\|^{k-t}. \quad (4)$$

Điều kiện vanishing và exploding gradients

Xét giới hạn khi độ dài chuỗi tăng ($T \rightarrow \infty$):

Nếu $\|W_{hh}\| < 1$:

$$\|W_{hh}\|^{k-t} \rightarrow 0 \quad \Rightarrow \quad \text{gradient vanishing.}$$

Nếu $\|W_{hh}\| > 1$:

$$\|W_{hh}\|^{k-t} \rightarrow \infty \quad \Rightarrow \quad \text{gradient exploding.}$$

Nếu $\|W_{hh}\| \approx 1$:

gradient được duy trì ổn định.

Trong thực hành, điều kiện này tương đương với bán kính phổ (spectral radius) của W_{hh} .

Tóm tắt vấn đề

Vanishing và exploding gradients là hệ quả trực tiếp của việc:

- Nhân lặp cùng một ma trận W_{hh} qua nhiều bước thời gian.

- Gradient là tổng của các chuỗi lũy thừa ma trận.

Do đó, việc huấn luyện RNN chuỗi dài gặp khó khăn ngay cả trong trường hợp tuyến tính đơn giản.

Giải pháp 1: Gradient Clipping

Gradient clipping giới hạn độ lớn gradient trong quá trình cập nhật.

Cụ thể:

$$g \leftarrow \frac{g}{\max\left(1, \frac{\|g\|}{\tau}\right)}, \quad (5)$$

trong đó g là gradient và τ là ngưỡng.

Đặc điểm:

- Không loại bỏ nguyên nhân gốc rễ.

- Ngăn gradient exploding trong thực tế.

- Dễ cài đặt, được dùng phổ biến.

Giải pháp 2: Long Short-Term Memory (LSTM)

LSTM thay đổi kiến trúc RNN bằng cách:

- Tạo đường truyền gradient gần như tuyến tính theo thời gian.

- Tránh nhân lặp trực tiếp bởi cùng một ma trận.

Trạng thái cell c_t được cập nhật dưới dạng:

$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t, \quad (6)$$

trong đó f_t là forget gate.

Gradient có thể đi xuyên qua nhiều bước thời gian mà không bị vanishing hoặc exploding nghiêm trọng.