

Recitation 3 Feb 23

Agenda

- Recurrent Neural Network
- BP in RNN
- LSTM

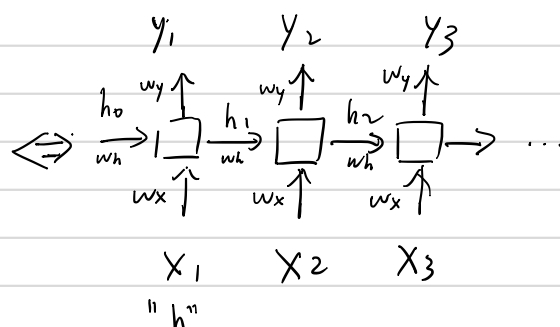
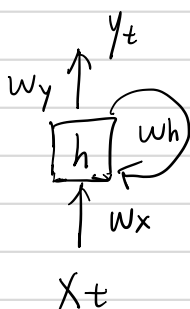
RNN

• General Structure

x_t : input at time t

y_t : output at time t

h_t : hidden state of RNN at time t



$$h_t = f_1(x_t, h_{t-1}; w_h, w_x)$$

$$y_t = f_2(x_t, h_{t-1}; w_y)$$

• RNN for char-level language model (CLLM)

$$\max_{w_h, w_y} p(x_1, x_2, x_3, \dots, x_t) \quad (\text{likelihood})$$

$$p(\text{"hello"}) \uparrow$$
$$p(\text{"htktt"}) \downarrow$$

$$p(x_1, x_2, \dots, x_t) = \prod_{i=1}^t p(x_i | x_{i-1}, \dots, x_1)$$

\Rightarrow • Use h_{t-1} to capture history input x_1, x_2, \dots, x_{t-1}

$$y_t = p(x_t | h_{t-1}) \approx p(x_t | x_{t-1}, \dots, x_1)$$

• Evaluate the loss by cross entropy $\sum_t E(y_t, \hat{y}_t)$

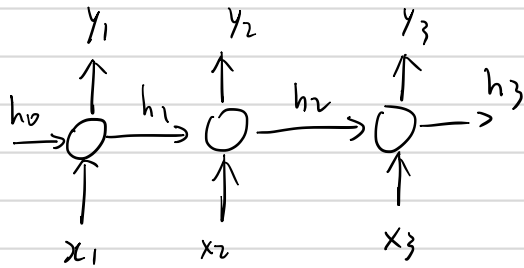
• Canonical RNN (for CLLM)

$$f_1 = \tanh(\cdot; w_h, w_x) \Rightarrow h_t = \tanh(W_{xh} x_t + W_{hh} h_{t-1})$$

$$f_2 = \text{softmax}(\cdot; w_y, b_y) \Rightarrow \underline{y_t = \text{softmax}(W_{yh} h_t + b_y)}$$

embedding

BP in RNN



$$\frac{\partial E}{\partial w_y}$$

$$\frac{\partial E}{\partial w_y} \stackrel{(1)}{=} \sum_t \frac{\partial E_t}{\partial w_y} \quad \frac{\partial E_t}{\partial y_t} \quad \frac{\partial y_t}{\partial w_y}$$

(2) (3) ✓ ✓

$$E = E_1 + E_2 + E_3 \quad (1)$$

$$E_t = E(y_t, \hat{y}_t) \quad (2)$$

$$y_t = \text{softmax}(W_y h_t + b_y) \quad (3)$$

$$h_t = \tanh(W_h h_{t-1} + W_x x_t) \quad (4)$$

$$\frac{\partial E}{\partial w_x}$$

$$\frac{\partial E}{\partial w_x} = \sum_t \frac{\partial E_t}{\partial w_x} \quad \frac{\partial E_t}{\partial y_t} \quad \frac{\partial y_t}{\partial h_t} \quad \frac{\partial h_t}{\partial w_x}$$

(2) (3) (4) ✓ ✓ ?

$$h_3 = f(x_3, h_2, W_h, W_x)$$

$$h_2 = f(x_2, h_1, W_h, W_x)$$

$$h_1 = f(x_1, h_0, W_h, W_x)$$

$$\Rightarrow \frac{\partial h_3}{\partial w_x} = \frac{\partial h_3}{\partial w_x} + \frac{\partial h_3}{\partial h_2} \frac{\partial h_2}{\partial w_x} + \frac{\partial h_3}{\partial h_2} \frac{\partial h_2}{\partial h_1} \frac{\partial h_1}{\partial w_x}$$

BP through time (BPTT) can vanish or explode

$$\frac{\partial E}{\partial w_h} \text{ similar}$$

• Fight against Vanishing / Exploding gradient

• Good activation (ReLU)

• Clip the gradient

• LSTM / GRU cell instead of tanh

• TBPTT (Truncated BPTT)

• Embedding

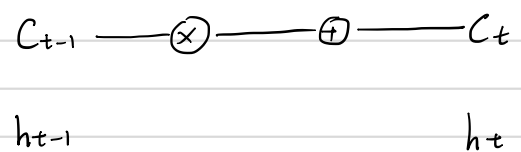
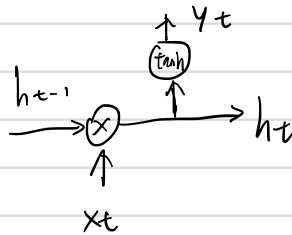
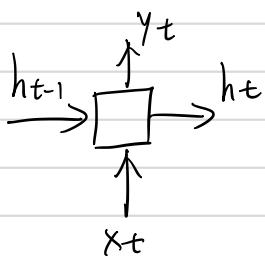
LSTM

Motivation: Learn longer dependencies (RNN in theory can, but doesn't in practice)

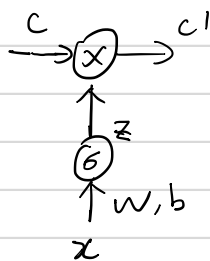
Hochreiter 1991 / Bengio 1994

eg 1: "Today is a good —"

eg 2: "I live in France. I speak —"



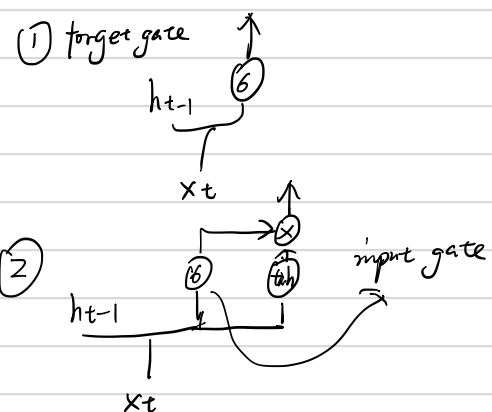
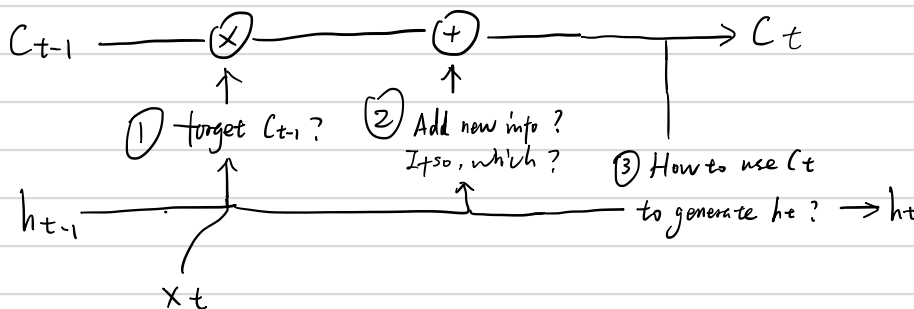
Gates



$$z = \sigma(WX + b) \in [0, 1]^{m \times n}$$

$C' = C \cdot z$ controls which part / how much of C flows through

Use "Gates" to control RNN when to forget / remember



Variants

- GRU \rightarrow ① combine forget and input gate ② merge h_t and C_t
- Many others.

