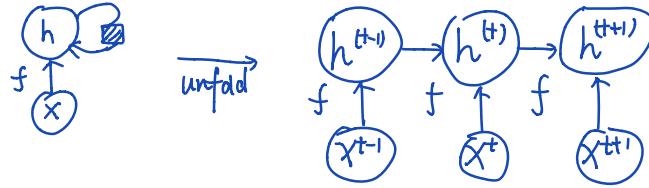


Unfold Computational Graphs

- Two representations of Recurrent Neural Networks



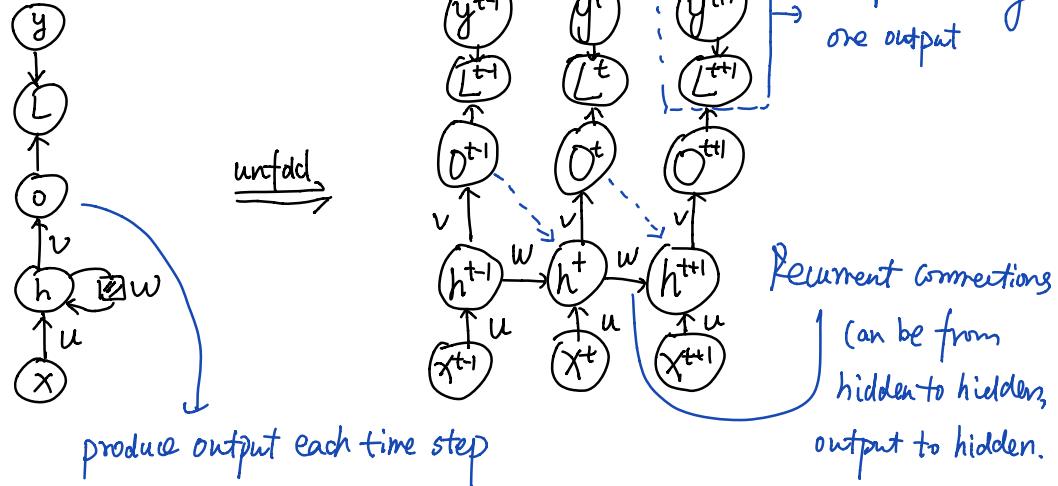
$$h^{(t)} = f(h^{(t-1)}, x^{(t)}, \theta)$$

The unfolding representation can clearly show the operations and shared weights, also can represent fixed length input and output.

- Here we can see that the operation at time step t , actually involves information from $1, \dots, t-1$, directly and indirectly.

Recurrent Neural Networks

Design patterns



- Running forward and computing the loss

$$\left. \begin{aligned} a^{(t)} &= b + w h^{(t-1)} + u x^{(t)} \\ h^{(t)} &= \tanh(a^{(t)}) \\ o^{(t)} &= c + v h^{(t)} \\ \hat{y}^{(t)} &= \text{softmax}(o^{(t)}) \end{aligned} \right\} \text{Forward Operation}$$

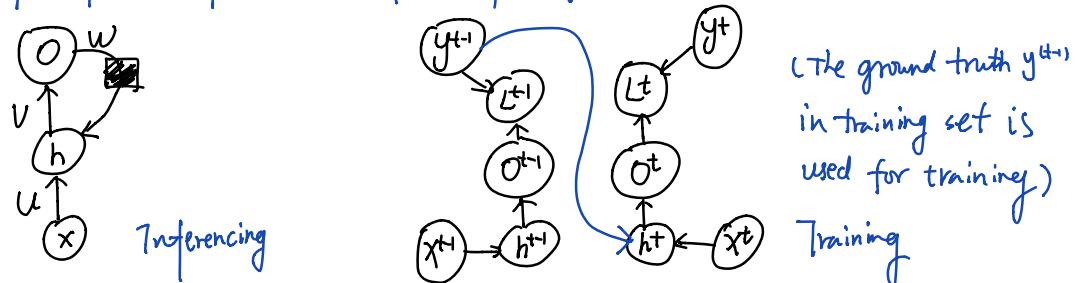
$$\begin{aligned} &L(\{x^{(1)}, \dots, x^{(T)}\}, \{y^{(1)}, \dots, y^{(T)}\}) \\ &= \sum_t l^{(t)} \\ &= - \sum_t \log P_{\text{model}}(y^{(t)} | x^{(1)}, \dots, x^{(t)}) \end{aligned}$$

Negative log likelihood ←

- Teacher Forcing and Networks with Output Recurrence.

Focus on training and analysis a specific RNN model structure.

The output of time step $t+1$ is one of the input of h^{t+1}



- Limits of teacher forcing

- Low ability of representation (h^t need to remember all the past)
- Low perform with training and testing difference is big.

- Computing Gradient of an RNN

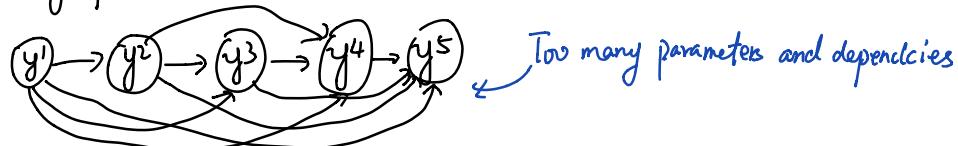
Computing the gradient of hidden units:

$$\nabla_{h^{t+1}} L = \underbrace{\left(\frac{\partial h^{t+1}}{\partial h^{t+1}} \right)^T (\nabla_{h^{t+1}} L)}_{\text{pass to } h^{t+1}, \text{ influencing } L^{t+1}} + \underbrace{\left(\frac{\partial o^{t+1}}{\partial h^{t+1}} \right)^T \nabla_{o^{t+1}} L}_{\text{self to } L^{t+1}}$$

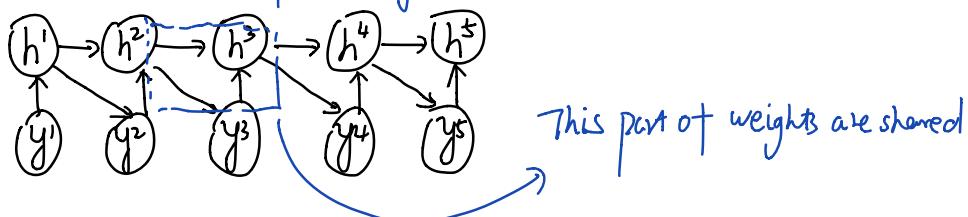
Other gradients on W, V, U, b, c , can be computed based on this.

- Recurrent Networks as Directed Graphical Models.

A directed graphical model with no recurrent hidden units:



↓ After adding recurrent hidden units



- Determining the output length
 - Definition a stop symbol
 - Each output associated with a binary variable indicating end or not.
 - Output an integer "t" of the length of sequence.
- Modeling Sequence Conditioned on Contexts with RNNs
 - Fixed length non-sequential input, output a sequence describing input, we input with an encoding/embedding style to the hidden units and recurrent use output and hidden units to next output (hidden timestep).
 - Sequential fixed length input and output, each hidden unit receives an input x , and output a y , y can also be as input to next-step hidden units for more representative information (not teacher forcing, cause we have hidden to hidden connection)

Bidirectional RNN

Adding backward connection from $g^{t+1} \rightarrow g^t$. namely considering another information flow from x^t to x^1 parallelly, each output is affected by both past and future.

Sequence to Sequence

- Encoder: An RNN inputting sequential x^t , and output a vector representation C , noticeable vector
- Decoder: As image caption, decoding the embedding vector C into sequential output y .
- C : often called context, attention mechanism to solve the drawbacks of fixed length context set up.

Deep Recurrent Networks

W, U, V (connection between hidden, output and input) are best designed as Neural Networks

Recursive Neural Networks

Tree-structured, recursive operation and finally reaching solely one output.

The challenge of Long-Term Dependencies

$$h^{(t)} = W h^{(t-1)}$$

$$\hookrightarrow h^{(t)} = [W^T]^T h^{(0)} = Q \Lambda^t Q h^{(0)}$$

where the exploding/vanishing occurs.

Echo State Networks

- Spectral Radius:

Jacobian J , $\lambda = \max(\omega_i)$, where ω is eigen-values of Jacobian.

When going backward or forward, the value of spectral radius of J , reveals the variation rate of updating steps.

- Echo State

Control spectral radius in a reasonable value when going backward.

Strategies for multiple time scales

- Adding skip connections / Removing connections then adding skip connections.
- Leaky units controlling how much information will preserve from the past.

$$U^{(t)} \leftarrow d U^{(t+1)} + (1-d) V^{(t)}$$

LSTM and other Gated RNN (GRU)

One principle: Adding gate between input to hidden, hidden to hidden and hidden to output.

$$\left. \begin{array}{l} f^{(t)} = \sigma(b^f + \sum_j U_j^f x_j + \sum_j W_j^f h^{(t-1)}) \\ S^{(t)} = f^{(t)} S^{(t-1)} + g^{(t)} \sigma(b + \sum_j U_j x_j + \sum_j W_j h^{(t-1)}) \end{array} \right\}$$

(Input/Output gate follow the same logic consisting x and h)

$$G^{(t)} = f^{(t)} S^{(t-1)} + g^{(t)} \sigma(b + \sum_j U_j x_j + \sum_j W_j h^{(t-1)})$$

GRU:

$$h^{(t)} = \underbrace{U^{(t-1)} h^{(t-1)}}_{\text{update gate}} + (1 - U^{(t-1)}) \underbrace{\sigma(b + \sum_j U_j x_j^{(t)} + \sum_j W_j Y_j^{(t)} h_j^{(t-1)})}_{\text{reset gate}}$$

Optimizing for gradient exploding / Cliff error surface

Second-order method (Newton method) can eliminate, but replaceable by other models with lower cost.

- Clipping gradients

$$\text{if } \|g\| > v, \text{ then } g \leftarrow \frac{gv}{\|g\|}$$

- Regularizing to Encourage Information Flow

$$\text{Adding regularizing term } \Omega = \sum_t \left(\frac{\|\nabla_{h^{(t)}} L \frac{\partial h^{(t)}}{\partial h^{(t+1)}}\|}{\|\nabla_{h^{(t)}} L\|} - 1 \right)^2$$

Explicit Memory

A task networks to control co-efficient of different sub-modules (memory cells).

Recurrent model based