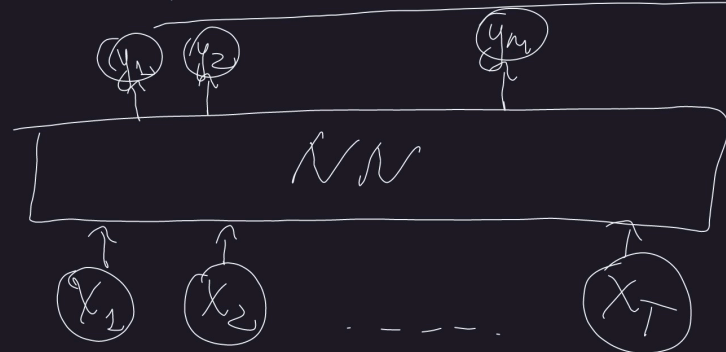


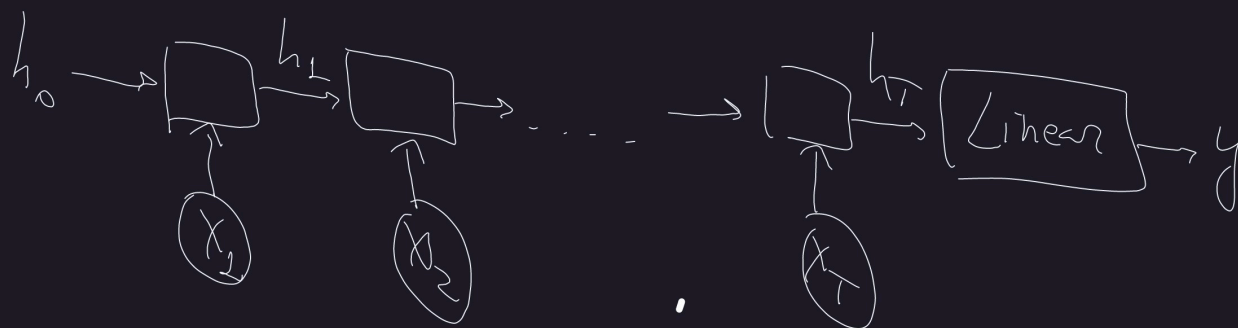
Recurrent Neural Networks and Attention



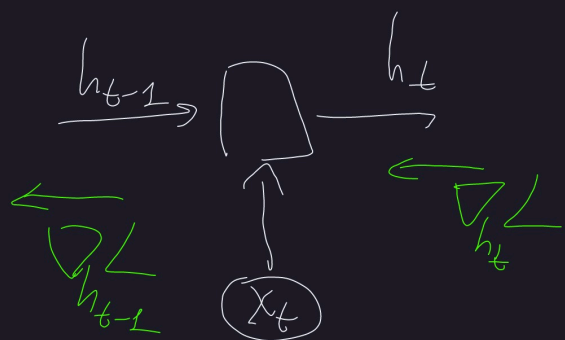
$$h_t = g(W_x x_t + W_h h_{t-1} + b)$$

$$x_t \in \mathbb{R}^{\text{input}}$$
$$h_t \in \mathbb{R}^{\text{hidden}}$$

$$W_x \in \mathbb{R}^{\text{hidden} \times \text{input}}$$
$$W_h \in \mathbb{R}^{\text{hidden} \times \text{hidden}}$$
$$b \in \mathbb{R}^{\text{hidden}}$$



RNN



$$\nabla_{h_{t-1}} L = W_h^T \left(\frac{\partial g}{\partial z_t} \right)^T \nabla_{h_t} L = W_h^T \frac{\partial g}{\partial z_t} \cdot W_h^T \frac{\partial g}{\partial z_{t-1}} \nabla_{h_{t-1}} L =$$

$$\nabla_{h_t} L^T = \left(\prod_{i=t}^T W_h^T \frac{\partial g}{\partial z_i} \right) \cdot \nabla_{h_T} L$$

$$dL = \nabla_{h_t} L^T dh_t = \nabla_{h_t} L^T \frac{\partial g}{\partial z_t} W_h \cdot dh_{t-1}$$

$$h_t = g(\underbrace{W_x x_t + W_h \cdot h_{t-1} + b}_{z_t})$$

$$\|\nabla_{h_{t-1}} L\| \leq \prod_{i=t}^T \|W_h\| \cdot \left\| \frac{\partial g}{\partial z_i} \right\| \cdot \|\nabla_{h_T} L\|$$

exploding gradients: gradient clipping

$$\text{clip}_{\epsilon}(g) = \begin{cases} \frac{g}{\|g\|} \cdot \epsilon, & \text{if } \|g\| > \epsilon \\ g, & \text{otherwise} \end{cases}$$

vanishing gradients: ReLU activ. +

orth. W_h : - init.

- regul. $+\lambda \|W_h^T W_h - I\|_F^2$

- reparam. $\expm(V - V^T)$

$$W^T W = I$$

LSTM

$$\frac{\partial c_t}{\partial c_{t-1}} = I$$

Long Short Term Memory Network

c_t - Long term memory

• h_t - short-term memory

$$c_t = f_t \odot c_{t-1} + g_t \odot i_t$$

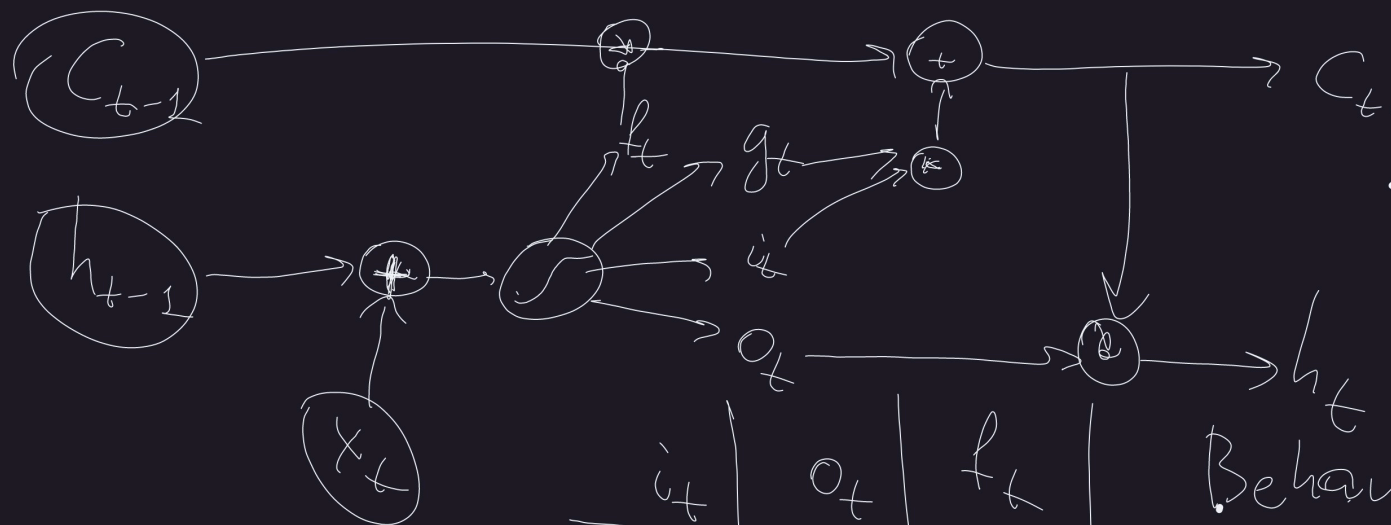
$$\cdot h_t = o_t \odot \tanh(c_t)$$

$$g_t = g(W_x^g x_t + W_h^g h_{t-1} + b^g)$$

$$f_t = \sigma(\dots)$$

$$i_t = \sigma(W_x^i x_t + W_h^i h_{t-1} + b^i) \in (0, 1) \quad \text{Init: } b^R \gg 1$$

$$o_t = \sigma(W_x^o x_t + W_h^o h_{t-1} + b^o)$$



i_t	o_t	f_t
+	+	-
+	-	+
-	+	+
-	-	+
-	-	-

Behavior of LSTM cell

RNN

Storing inform. to memory
 Loading inform.
 preserving inf.
 erasing inf.

GRU

gated recurrent unit

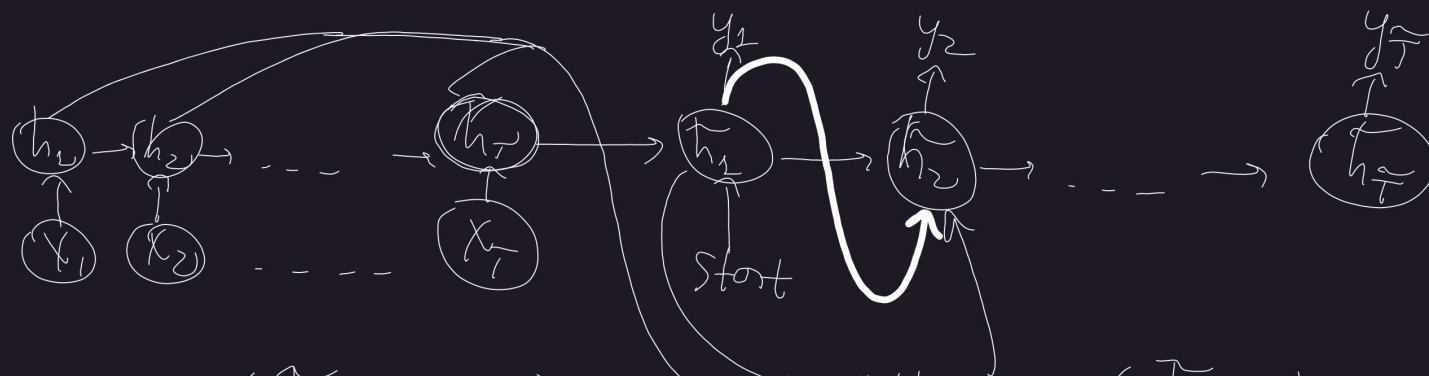
$$z_t = \sigma(W_x^z \cdot x_t + W_h^z \cdot h_{t-1} + b^z)$$

$$u_t = \sigma(\dots)$$

$$g_t = g(W_x^g \cdot x_t + W_h^g (h_{t-1} \odot z_t) + b^g)$$

$$h_t = (1 - u_t) \odot g_t + u_t \odot h_{t-1} \quad \text{Init: } b^u \gg 1$$

a b a e e b c d → a b a e e b c d a b a e e b c d



$$\text{score}(\tau_{t-1}, h_i) \in \mathbb{R} \Rightarrow \text{Attention}(\tau_{t-1}, h_1, \dots, h_T)$$

$$d_i = \frac{\exp(\text{score}(\tau_{t-1}, h_i))}{\sum_{j=1}^T \exp(\text{score}(\tau_{t-1}, h_j))}, i=1, \dots, T$$

$$\text{Att} = \sum_{i=1}^T d_i h_i$$

$$\text{score}(x, y): \quad (1) \quad x^T y$$

$$(3) w^T \tanh(W_x x + W_y y) + b \quad (2) \quad \frac{x^T y}{\sqrt{\dim(x)}}$$

$$x^T y = \sum_{i=1}^d x_i y_i$$

$$x_i, y_i \sim \mathcal{N}(0, 1)$$

$$\begin{aligned} \text{Var}(x^T y) &= \text{Var} \sum_i x_i y_i = \mathbb{E} x^T y = \mathbb{E} \sum_i x_i y_i = \\ &= \sum_{i=1}^d \text{Var}(x_i) \cdot \text{Var}(y_i) = \sum_i \mathbb{E} x_i \mathbb{E} y_i = 0 \end{aligned}$$

$$\text{Var} \left(\frac{x^T y}{\sqrt{\dim(x)}} \right) = 1 \quad \text{I} = d$$

