



similar to the example of a RNN, gradients are backpropagated through time.

dwx , db , and dwx are internal state variables shared across the LSTM cells in time, so the gradients get summed

$$\therefore dwx += dwx_{-t} (D, 4H)$$

$$dwx += dwx_{-t} (H, 4H)$$

$$db += db_{-t} (4H,)$$

for is not enforced by the close 1 loss and

dx is just returned from the step backprop and stored in a variable. We don't accumulate them because

each time step input is unique.

$$\therefore dx_{(N,H)}(t;T,:) = dx_{(N,T,D)} \rightarrow (N,D)$$

- $dx_{prev,h}$ is returned from the step backprop and accumulated. dx_0 is the accumulation at the first timestep.