



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

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

$$\therefore dwx_t = dwx_{t-1}$$

$$dwh_t = dwh_{t-1}$$

$$db_t = db_{t-1}$$

but is not enforced for the cell 1, because

$ax$  is just returned from the step backprop and stored in a variable. We don't accumulate them because each time step input is unique.

- $dprer.h$  is returned from the step backprop and accumulated.  $dho$  is the accumulation at the first timestep.