

Vanilla RNN

$$h_t = \tanh\left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right)$$

Vanilla LSTM

$$\begin{pmatrix} hi \\ hf \\ ho \\ hc \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh(\cdot) \end{pmatrix} \begin{pmatrix} W_i \\ W_f \\ W_o \\ W_c \end{pmatrix} \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = hf \odot c_{t-1} + hi \odot hc$$

$$h_t = ho \odot \tanh(c_t)$$

c_t , cell state, hidden from the outside world

hf , forget gate, indicator to keep or forget previous cell state element

hi , input gate, indicator to edit cell state element

hc , cell gate, increment or decrement cell state element by a value between -1 and 1

ho , output gate, indicator to reveal previous/edited cell state element to outside world

$h_t = (1, H) = \text{number of LSTM layer neurons}$

$x_t = (1, D) = \text{number of items in vocabulary}$

$\hat{y}_t = (1, D) = \text{output is probabilities over vocabulary set}$

$Z = H + D = \text{concatenated size}$

$W_f, W_i, W_o, W_c = (Z, H) \quad b_f, b_i, b_o, b_c = (1, H) \quad hf_t, hi_t, ho_t, hc_t = (1, H)$

$W_v = (H, D) \quad b_v = (1, D)$

$$L_k = - \sum_{t=k}^T \sum_j y_{t,j} \log \hat{y}_{t,j}$$

$$L = L_1$$

$z_t = [h_{t-1}, x_t]$ $hf_t = \sigma(z_t \cdot W_f + b_f)$ $hi_t = \sigma(z_t \cdot W_i + b_i)$ $ho_t = \sigma(z_t \cdot W_o + b_o)$ $hc_t = \tanh(z_t \cdot W_c + b_c)$ $c_t = hf_t \odot c_{t-1} + hi_t \odot hc_t$ $h_t = ho_t \odot \tanh(c_t)$ $v_t = h_t \cdot W_v + b_v$ $\hat{y}_t = \text{softmax}(v_t)$	$dv_t = \hat{y}_t - y_t$ $dh_t = dv_t \cdot W_v^T + dh'_t$ $dho_t = dh_t \odot \tanh(c_t)$ $dc_t = dh_t \odot o_t \odot (1 - \tanh^2(c_t)) + dc'_t$ $dhf_t = dc_t \odot c_{t-1}$ $dhi_t = dc_t \odot hc_t$ $dhc_t = dc_t \odot hi_t$ $dhf'_t = hf_t \odot (1 - hf_t) \odot dhf_t$ $dhi'_t = hi_t \odot (1 - hi_t) \odot dhi_t$ $dhc'_t = (1 - hc_t^2) \odot dhc_t$ $dho'_t = ho_t \odot (1 - ho_t) \odot dho_t$ $dz_t = dhf'_t \cdot W_f^T + dhi'_t \cdot W_i^T + dho'_t \cdot W_o^T + dhc'_t \cdot W_c^T$ $dW_v = h_t^T \cdot dv_t \quad db_v = dv_t$ $dW_f = z^T \cdot dhf'_t \quad db_f = dhf'_t$ $dW_i = z^T \cdot dhi'_t \quad db_i = dhi'_t$ $dW_c = z^T \cdot dhc'_t \quad db_c = dhc'_t$ $dW_o = z^T \cdot dho'_t \quad db_o = dho'_t$ $[dh'_{t-1}, dx_t] = dz_t$ $dc'_t = hf_t \odot dc_t$
--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

yhat = softmax(v)

dv = yhat.copy()

dv[1, y_train] -= 1

$v = h @ W_v + b_v$	$dh = dv @ W_v.T \#(1,H)=(1,D)(D,H)$ $dh += dh_next$ $dW_v = h.T @ dv \#(H,D)=(H,1)(1,D)$ $db_v = dv * 1$
$h = h_o * \tanh(c)$	$dho = dh * \tanh(c)$ $dc = dh * h_o * dtanh(c)$ $dc += dc_next$
$c = h_f * c_prev + h_i * h_c$	$dhf = dc * c_prev$ $dhi = dc * h_c$ $dhc = dc * h_i$
$h_f = \text{sigmoid}(z @ W_f + b_f)$	$dhf *= \text{dsigmoid}(h_f)$ $dz = dhf @ W_f.T \#(1,Z)=(1,H)(H,Z)$ $dW_f = z.T @ dhf \#(Z,H)=(Z,1)(1,H)$ $db_f = dhf * 1 \#(1,H)=(1,H)$
$h_i = \text{sigmoid}(z @ W_i + b_i)$	$dhi *= \text{dsigmoid}(h_i)$ $dz += dhi @ W_i.T \#(1,Z)=(1,H)(H,Z)$ $dW_i = z.T @ dhi \#(Z,H)=(Z,1)(1,H)$ $db_i = dhi * 1 \#(1,H)=(1,H)$
$h_o = \text{sigmoid}(z @ W_o + b_o)$	$dho *= \text{dsigmoid}(h_o)$ $dz += dho @ W_o.T \#(1,Z)=(1,H)(H,Z)$ $dW_o = z.T @ dho \#(Z,H)=(Z,1)(1,H)$ $db_o = dho * 1 \#(1,H)=(1,H)$
$h_c = \tanh(z @ W_c + b_c)$	$dhc *= dtanh(h_c)$ $dz += dhc @ W_c.T \#(1,Z)=(1,H)(H,Z)$ $dW_c = z.T @ dhc \#(Z,H)=(Z,1)(1,H)$ $db_c = dhc * 1 \#(1,H)=(1,H)$ $dh_next = dz[:, 1:H] \#Z=H+D$ $dc_next = hf * dc$