hidden layers

Input  W

X

$X_2$

1
Bias

ReLu
b

tanh
b

~~Sigmoid~~
b

Sigmoid

Alt last layer

softmax

$\begin{bmatrix} - \\ - \end{bmatrix}$

last layer

2-layer

$$h = \partial \left( \underset{3 \times 2}{W} \cdot \underset{2 \times 1}{x} + b \right)$$
$$\underset{3 \times 1}{}$$

$$\hat{y} = Sigmoid(h)$$

$$ReLu(W_{r1} \cdot x_1 + W_{r2} x_2 + b)$$

# Language Modeling
## Predict the next word

Output $P(w|c) =$ [ $P(v_1)$ --------- $P(v_n)$ ] $\leftarrow$ Softmax $1 \times |v|$

$U: |v| \times d_h$

ReLU $1 \times d_h$

hidden layer [ $h_1$ ----- $h_n$ ]

$E_{mb}$ $Dd$

$d_h \times 2|v|$

Input trigram  W2U

[ $W_1$ ] $+$ [ $W_2$ ]

one

$1 \times 2 \cdot |v|$ [ next word ]
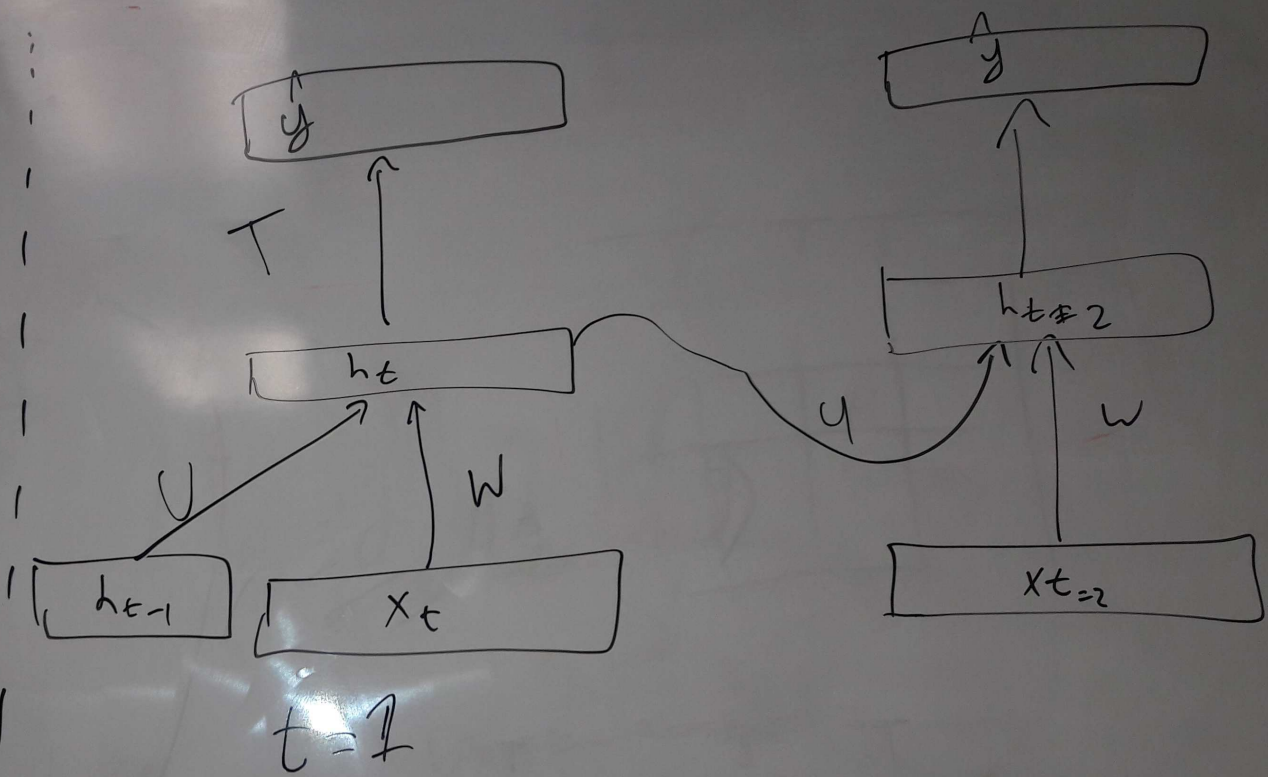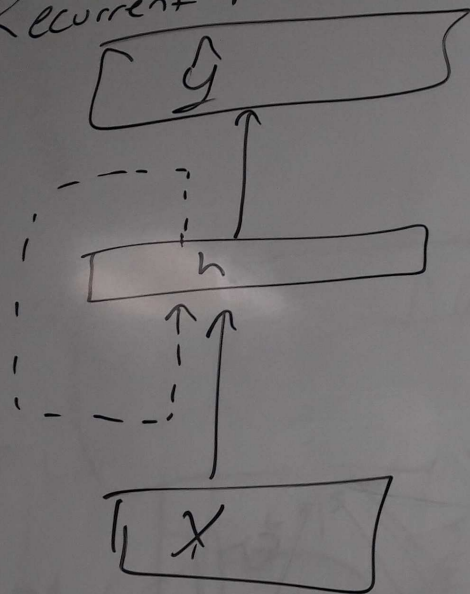
training

Loss     Cross-Entropy     $L_{ce}(y, \hat{y}) = -\log \dfrac{e^{z_i}}{\sum\limits_{i=1}^{4} e^{z_i}}$  negative
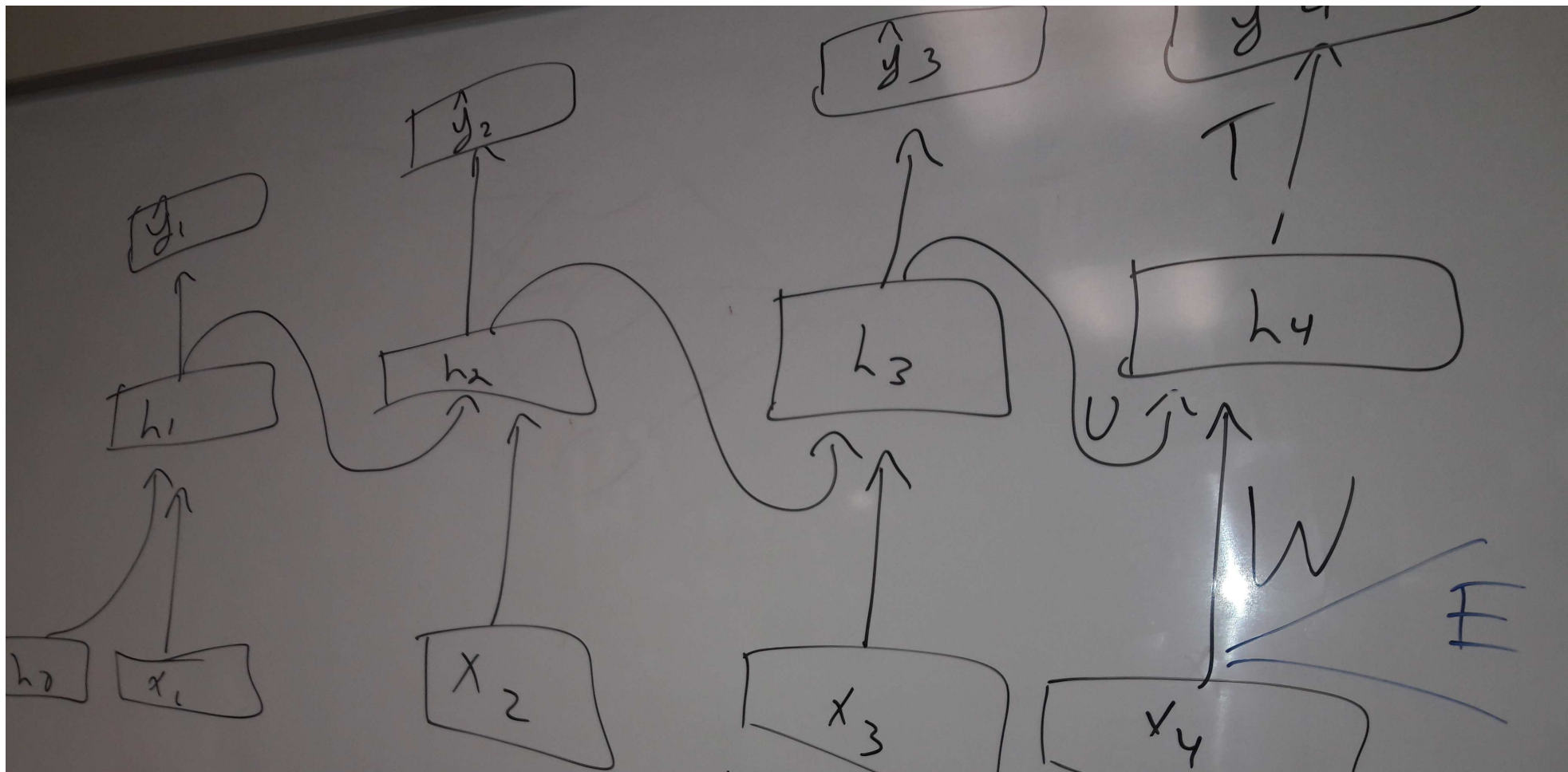log of softmax

Pytorch / torch
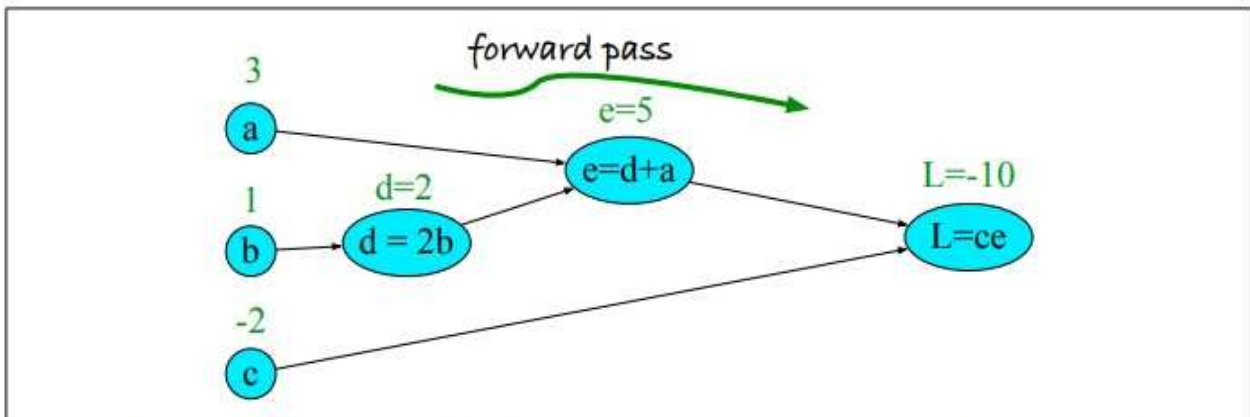(Python)    (Lua)

Tensorflow

Recurrent NN



$\hat{y}$

$h$

$X$

$\hat{y}$

$T$

$h_t$

$U$          $W$

$h_{t-1}$          $x_t$

$t=1$

$\hat{y}$

$h_{t=2}$

$U$          $W$

$x_{t=2}$

$W_1$  $W_2$  $W_3$
$t_{-y_1}$   $t_{-y_2}$  $t_{-y_1}$

$\hat{y}_1$

$\hat{y}_2$

$\hat{y}_3$

$\hat{y}_4$

$h_0$ $h_1$ $h_2$ $h_3$ $h_4$
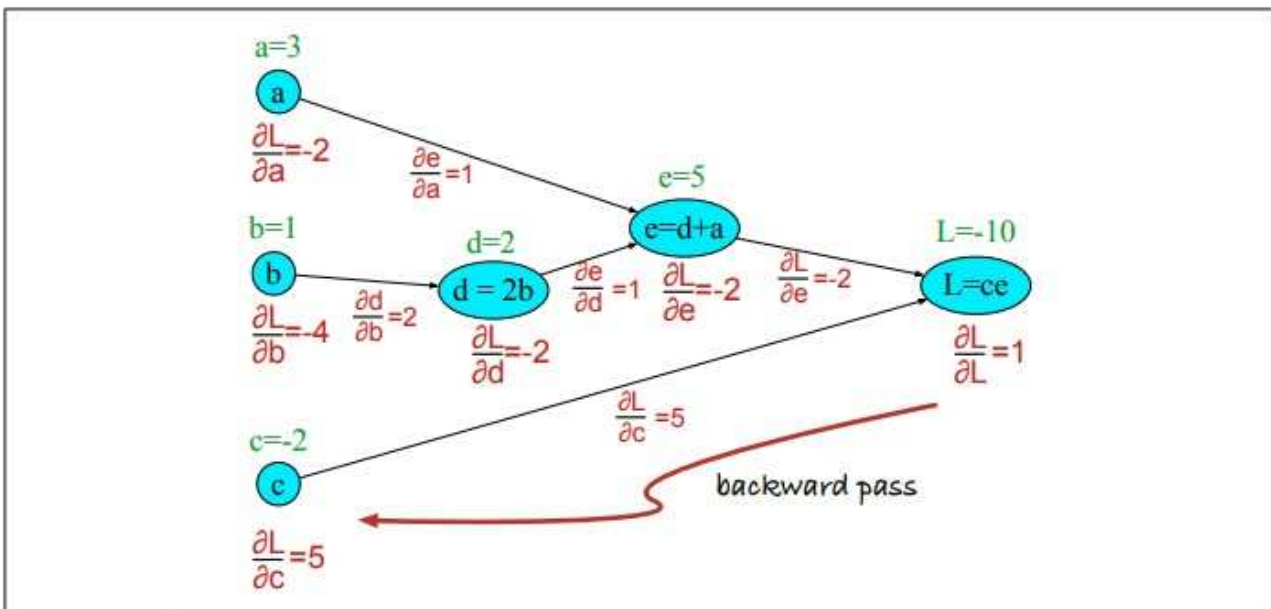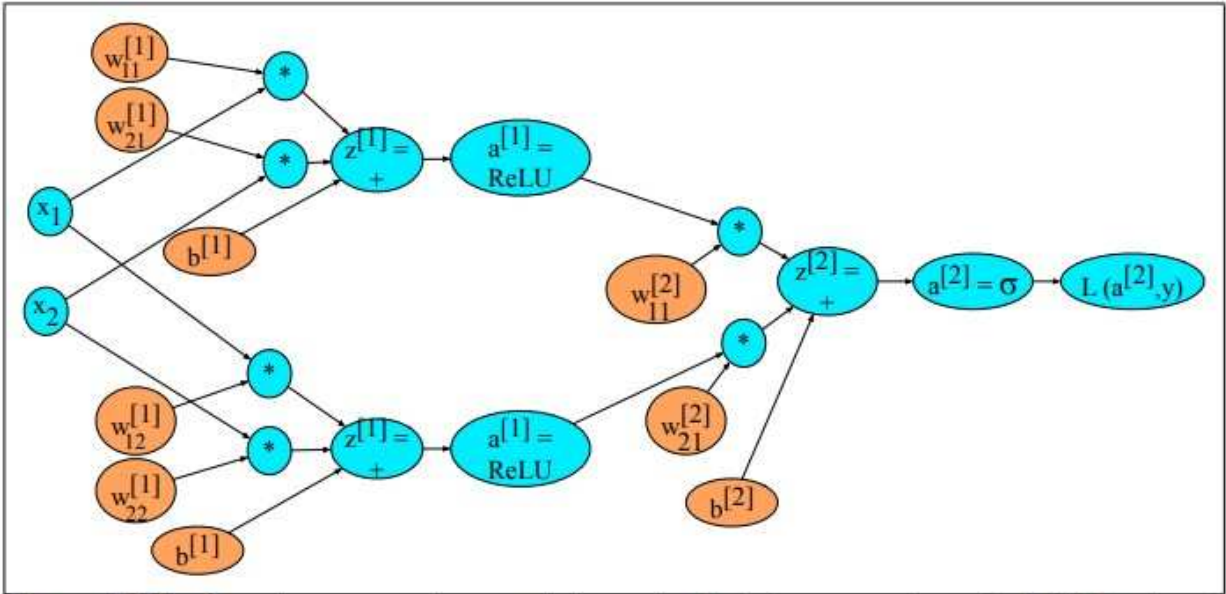
$x_1$ $x_2$ $x_3$ $x_4$

$T$

$U$

$W$

$E$

adjustable + raining properties

length —

Betching

**Figure 7.9** Computation graph for the function $L(a,b,c) = c(a+2b)$, with values for input nodes $a = 3$, $b = 1$, $c = -2$, showing the forward pass computation of $L$.



**Figure 7.10** Computation graph for the function $L(a,b,c) = c(a+2b)$, showing the backward pass computation of $\frac{\partial L}{\partial a}$, $\frac{\partial L}{\partial b}$, and $\frac{\partial L}{\partial c}$.

**Figure 7.11** Sample computation graph for a simple 2-layer neural net (= 1 hidden layer) with two input dimensions and 2 hidden dimensions.