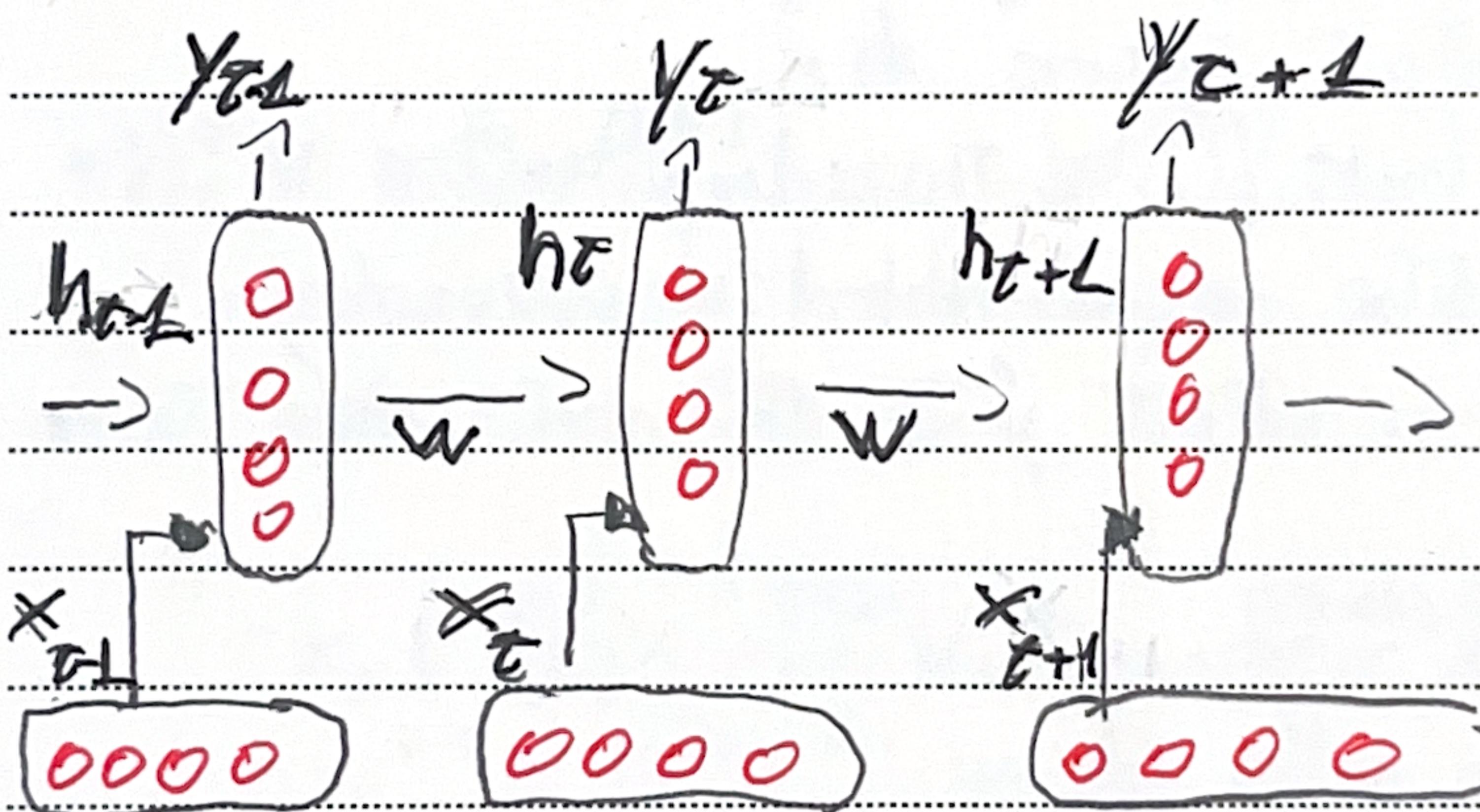


~~Topic 3: Recurrent Neural Networks~~

Recurrent Neural Network (RNN)

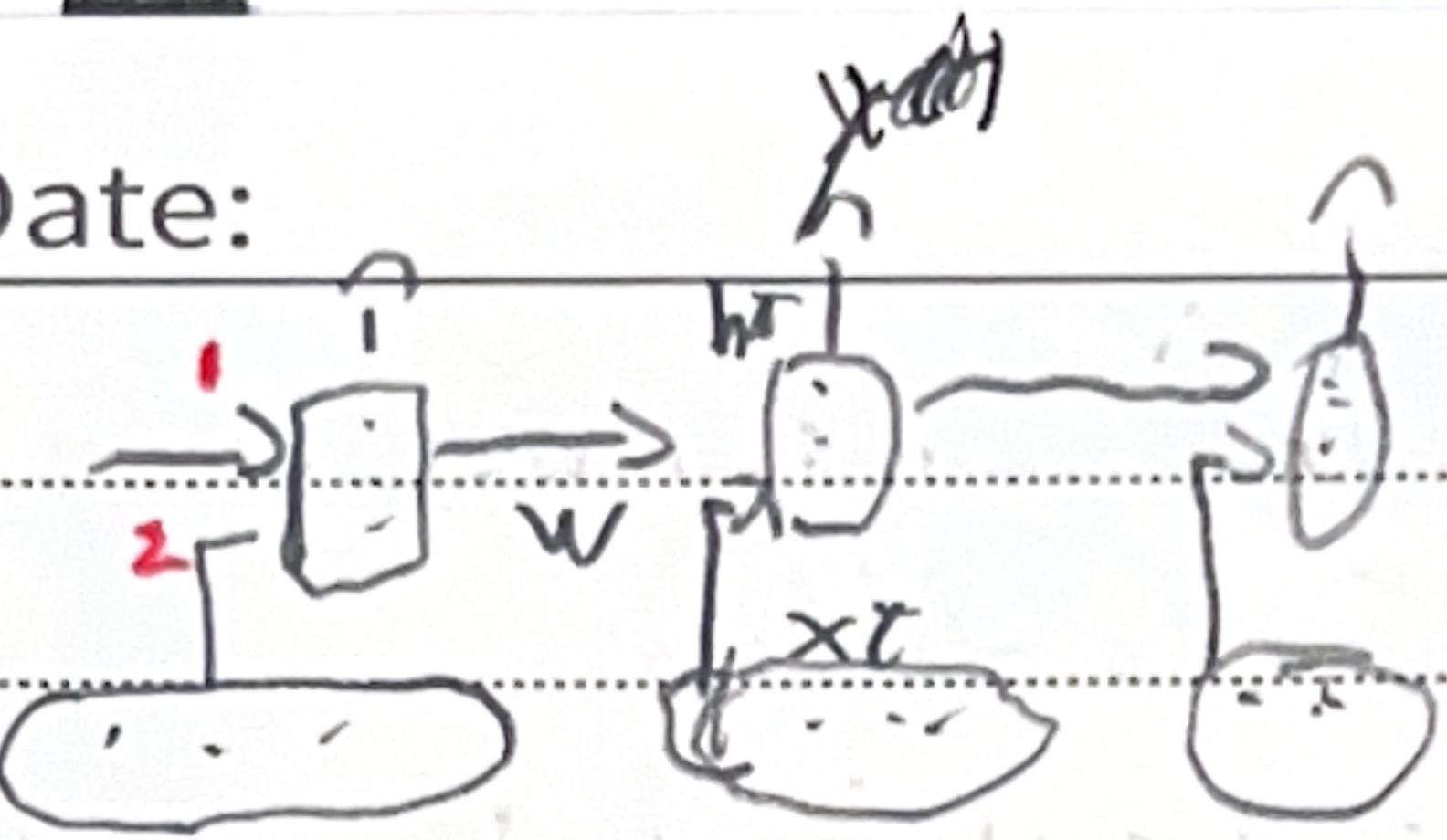
- Special Type of NN adapted to work for Time Series data or data that works in sequences.
- Capable of conditioning the model on all previous inputs.

A RNN with Three Time Steps



- Each vertical rectangle box is a hidden layer at time step t .
- Each layer holds a number of neurons, each of which performs a linear matrix operation on its inputs followed by a non-linear operation (e.g. Tanh or ReLU).
- At each Time step there are 2 inputs.

Date: _____



• AT each time step, we have 2 inputs

To the hidden layer

1 - the output of the previous layer h_{t-1}

2 - The input at that timestep x_t

3 - The former input is multiplied by a weight matrix $W^{(hh)}$

- And the latter by a weight matrix $W^{(hx)}$ to produce output features h_t

4 - These are multiplied with a weight matrix $W^{(s)}$

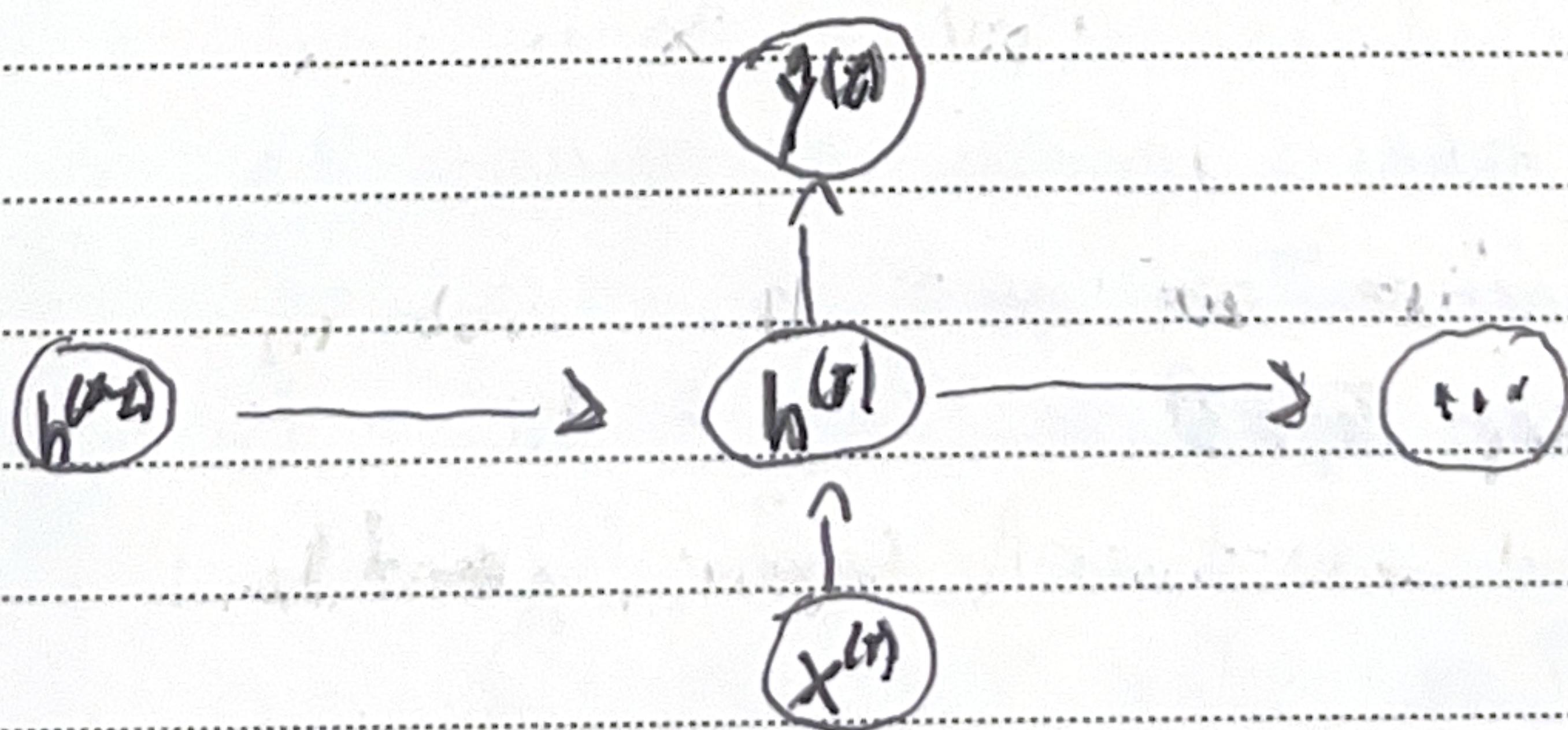
and through a softmax over the vocabulary

This produces output y_t of the next word.

$$3. h_t = \sigma(W^{(hh)} h_{t-1} + W^{(hx)} x_t)$$

$$4. y_t = \text{softmax}(W^{(s)} h_t)$$

• Inputs and outputs of each neuron illustrated here:



The same weights (W^h) and W^{chx} are applied repeatedly at each time step:

- number of parameters the model has to learn is less.
- independent of the length of the input sequence.
- i.e. deferring the curse of dimensionality.

Details of each parameter in the network

$X_1, \dots, X_{t-1}, X_t, X_{t+1}, \dots, X_T$: The word vectors corresponding to a corpus with T words

$h_t = \sigma(W^{hh} h_{t-1} + W^{chx} X_t)$: The relationship to compute the hidden layer output features at each time step t

$- X_t \in \mathbb{R}^d$: input word vector at time t .

$- W^{chx} \in \mathbb{R}^{D_h \times d}$: weights matrix used to condition the input word vector, X_t

$- W^{hh} \in \mathbb{R}^{D_h \times D_h}$: weights matrix used to condition the output of the previous time-step, h_{t-1}

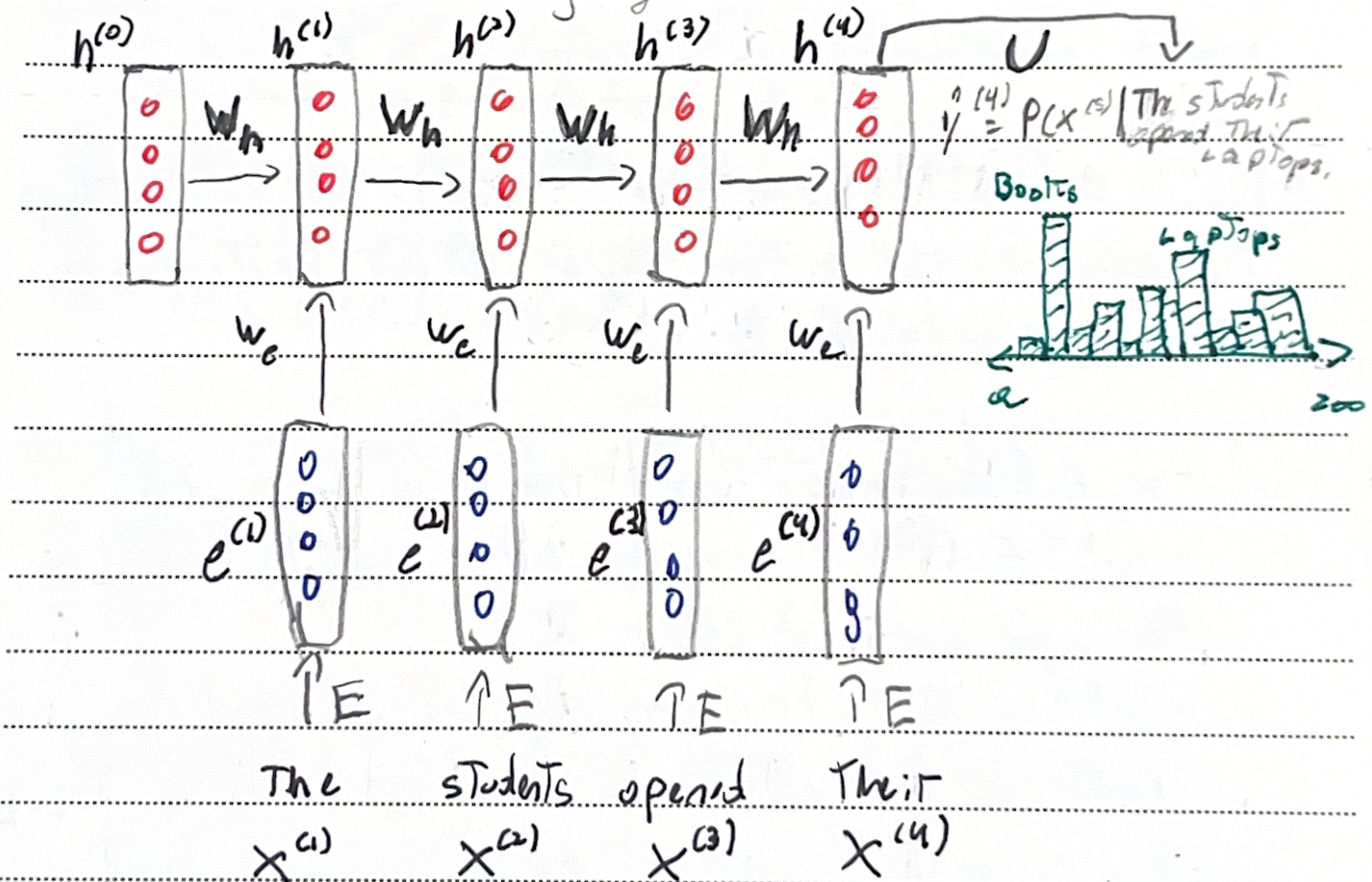
$- h_{t-1} \in \mathbb{R}^{D_h}$: output of the non-linear function at the previous time-step, $T-1$. $h_0 \in \mathbb{R}^{D_h}$ is an initialized vector for the hidden layer at time-step $T=0$.

$- \sigma()$: The non-linear function (e.g. sigmoid).

Date:

$\hat{y}_t = \text{Softmax}(W^{(s)} h_t)$: The output probability distribution over the vocabulary at each time-step t . Essentially, \hat{y}_t is the next predicted word given the document context score so far ($h_{0:T}$) and the last observed word vector $x^{(t)}$. Here, $W^{(s)} \in \mathbb{R}^{|V| \times D_h}$ and $\hat{y} \in \mathbb{R}^{|V|}$ where $|V|$ is the vocabulary.

Example of RNN Language Model



Here, the equivalent of $W^{(hh)}$ is W_h , $W^{(ht)}$ is U , and $W^{(s)}$ is W_e . E converts word inputs $x^{(t)}$ to word embeddings $e^{(t)}$. The final softmax over the vocabulary shows us the probability of various options for token $x^{(s)}$, conditioned on all previous tokens. The input could be much larger than 4-5 tokens.

RNN loss and perplexity

The loss function used in RNNs is often the cross entropy error.

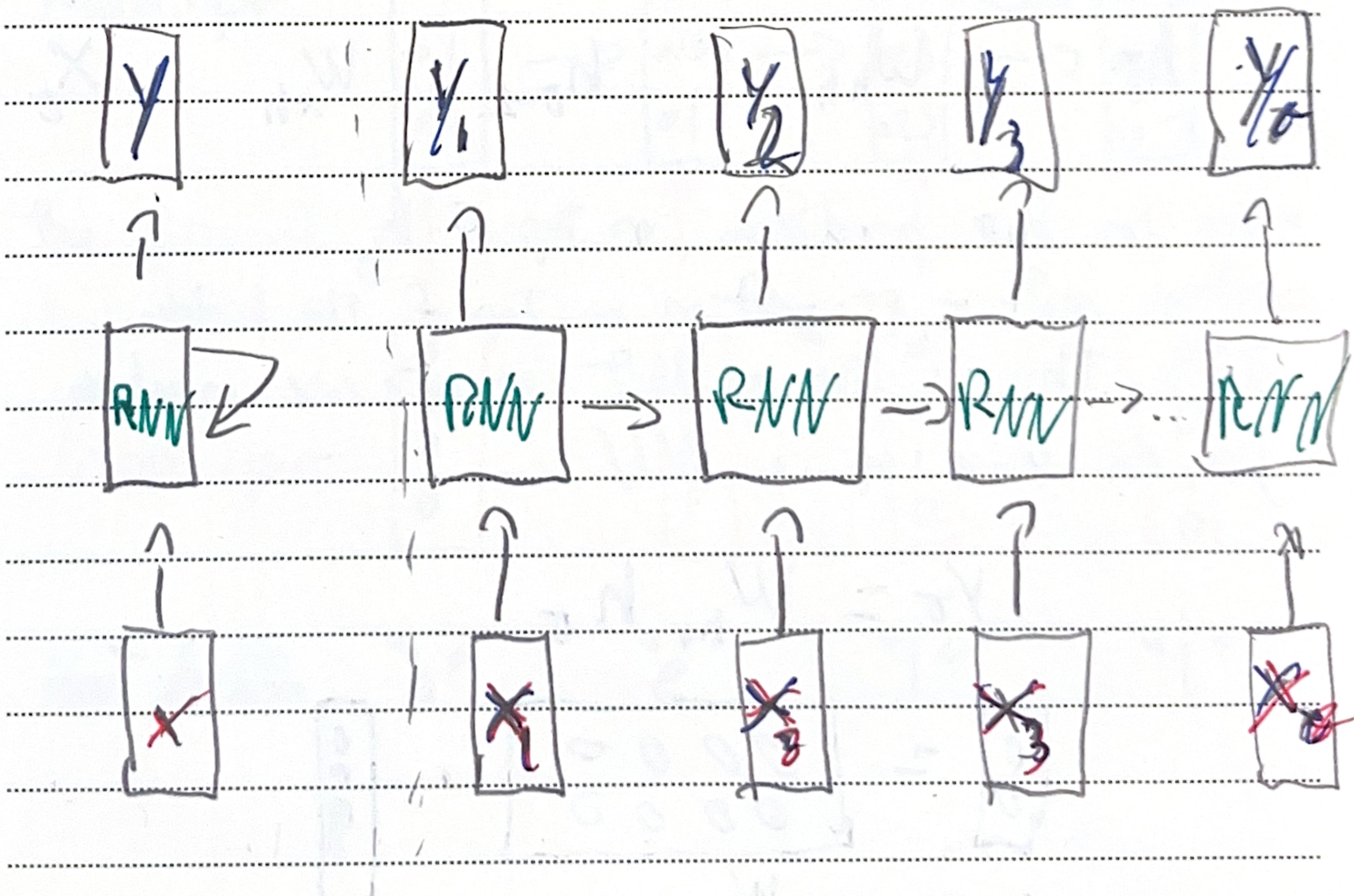
Sum of entire vocabulary at time step t

$$J(\theta) = - \sum_{j=1}^{|V|} Y_{t,j} \times \log (\hat{Y}_{t,j})$$

END TO END Example

Simplified

unrolled RNN



Date: Example RNN high level review.

$$h_t = f_w(h_{t-1}, x_t)$$

recurrent formula

- Simplest RNN (vanilla RNN), single hidden state h where we use a recurrent formula. That tells us how we should update our hidden state h .

$$h_t = \tanh(W_{hh} h_{t-1} + W_{xh} x_t)$$

$$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} = \tanh \left(\begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \cdot \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \cdot \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \right)$$

$$h_t \quad W_{hh} \quad h_{t-1} \quad W_{xh} \quad x_t$$

- We can base prediction on top of h_t by using just another matrix projection on top of the hidden state. This is the simplest complete case in which you can wire up a NN.

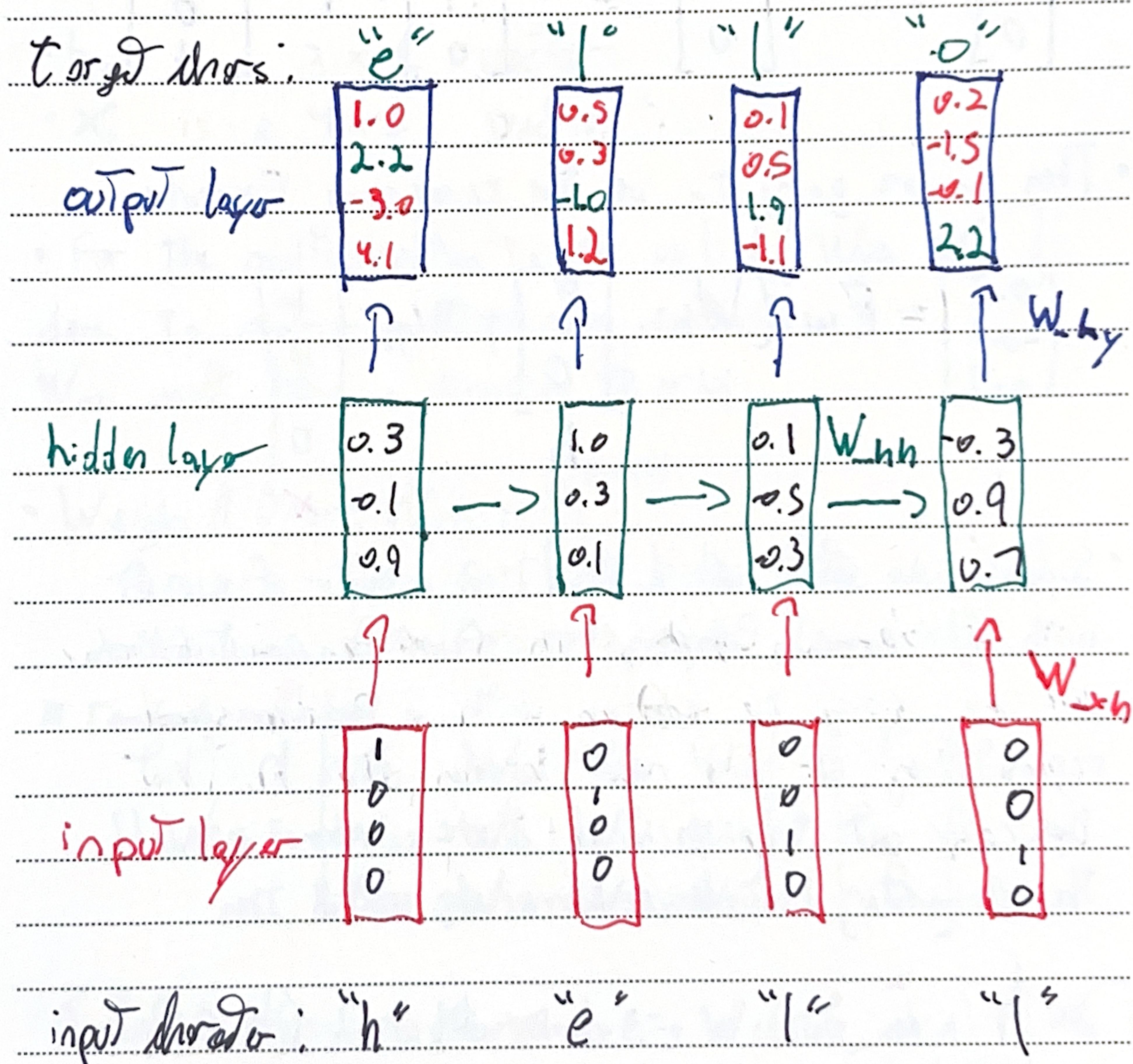
$$y_t = W_{hy} h_t$$

$$\begin{bmatrix} 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \cdot \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

Page _____

RNN example: character-level language model

- A Training sequence of just one string "hello" and we have a vocabulary $V \in \{h, e, l, o\}$ of 4 characters in the entire dataset.
- We are going to try to get an RNN to learn to predict the next character in the sequence of this training data.



Date:

- As shown in the last image, we will feed one character at a time into an RNN. First "h", then "e", then "l", and finally "o". All characters are encoded in the representation of a one-hot vector, where only one unique bit of the vector is turned on for each character.

$$\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} = "h" \quad \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = "e" \quad \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} = "l" \quad \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} = "o"$$

- Then we are going to use the recurrence formula:

$$\begin{bmatrix} 0.3 \\ -0.1 \\ 0.9 \end{bmatrix} = F_W \left(W_{hh} \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} + W_{xh} \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \right)$$

X

- Suppose we start off with h as a vector of size 3 with all zeros. By using the fixed recurrent formula, we are going to end up with a 3-dimensional representation of the next hidden state h . That basically at any point in time, summarizes all the characters that have come up until then.

- h (hidden state) is a 3-dimensional vector ($h \in \mathbb{R}^3$)

- The input vector ~~X~~ is a 4-dimensional vector ($x \in \mathbb{R}^4$)

The recursive formula is:

$$h_1 = \tanh(W_{hh} h_0 + W_{xh} x_1)$$

Since $h_0 = 0$, we simplify to:

$$h_1 = \tanh(W_{xh} \cdot x_1)$$

- h_1 is a 3×1 vector

- x_1 is a 4×1 vector

- For the multiplication to be valid, $[W_{xh} \cdot x_1]$

denote matrix multiplication constraints.

W_{xh} must have dimensions 3×4

- Weight Matrix; $W_{xh} \in \mathbb{R}^{3 \times 4}$

- 3 rows (one for each hidden state dimension)

- 4 columns (one for each input dimension)

* Randomly initialized

$$W_{xh} = \begin{bmatrix} W_{11} & W_{12} & W_{13} & W_{14} \\ W_{21} & W_{22} & W_{23} & W_{24} \\ W_{31} & W_{32} & W_{33} & W_{34} \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\text{Row 1} = W_{11} \cdot 1 + W_{12} \cdot 0 \dots = W_{11}$$

$$\text{Row 2} = W_{21} \cdot 1 + W_{22} \cdot 0 \dots = W_{21}$$

$$\text{Row 3} = W_{31} \cdot 1 + W_{32} \cdot 0 \dots = W_{31}$$

$$W_{xh} \cdot x = \begin{bmatrix} W_{11} \\ W_{21} \\ W_{31} \end{bmatrix}$$

Date:

Assume

$$h_1 = \tanh \begin{bmatrix} w_{11} \\ w_{21} \\ w_{31} \end{bmatrix}$$

$$w_{11} = 0.3$$

$$w_{21} = -0.1$$

$$w_{31} = 1.5$$

- We randomly initialize weights only during first iter.

- Then we apply our activation function

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$1. \tanh(0.3) \rightarrow$$

$$2. \tanh(-0.1) \quad ① = \frac{e^{0.3} - e^{-0.3}}{e^{0.3} + e^{-0.3}}$$

$$3. \tanh(1.5) \quad ② = \frac{e^{1.5} - e^{-1.5}}{e^{1.5} + e^{-1.5}}$$

$$1. \tanh(0.3) = 0.3$$

$$\begin{bmatrix} 0.3 \\ -0.1 \\ 0.9 \end{bmatrix}$$

$$2. \tanh(-0.1) = -0.1 \rightarrow \begin{bmatrix} 0.3 \\ -0.1 \\ 0.9 \end{bmatrix} = h_1$$

$$3. \tanh(1.5) = 0.9$$

Time Step 2

$$h_2 = \tanh(W_{hh} \begin{bmatrix} 0.3 \\ -0.1 \\ 0.9 \end{bmatrix} + W_{xh} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix})$$

- W_{hh} is a 3×3 matrix x that is also

randomly initialized at step 1 in this example

$$W_{hh} \text{ is all } 0's = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

Summary of Forward pass (Time Stamps)

$$h_1 = \tanh \left(\frac{0.3}{0.1} \cdot \frac{0.07}{0.01} + \frac{0.3, 3.0, 0.1, -0.3}{-0.1, 0.3, -0.5, 1.5} \right)$$

Рекурсия Геджи.

$$h_t = \tanh(W_{hh} h_{t-1} + W_{xh} \cdot x_t)$$

$$h_2 = \tanh \left(W_{hh} \begin{bmatrix} 0.3 \\ -0.1 \\ 0.9 \end{bmatrix} + W_x h \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \right) = \begin{bmatrix} 0.3 \\ 0.1 \\ 0.1 \end{bmatrix}$$

$$h_3 = \tanh(w_{hn} \cdot \begin{bmatrix} 1.0 \\ 0.3 \\ 0.1 \end{bmatrix} + w_{xh} \cdot \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} - h_2 \cdot \begin{bmatrix} 0.1 \\ -0.5 \\ -0.3 \end{bmatrix})$$

$$h_4 = \tanh \left(W_{hh} \begin{bmatrix} 0.1 \\ -0.5 \\ -0.3 \end{bmatrix} + W_{xh} \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \right) = \begin{bmatrix} 0.9 \\ 0.7 \end{bmatrix}$$

Date:

- After computing the hidden states (h_t) using our recursive relation, we can compute y_t

- Step 2: Compute output logits (y_t)

* Matrix multiplication

$$y_t = W_{hy} h_t$$

- This gives logits for predicting the next character.

- We can do this for all logits at once

- We use the hidden-to-output weight matrix W_{hy} 4×3
 W_{hy} is a learnable parameter; it is first randomly initialized

Assume,

$$W_{hy} = \begin{bmatrix} 0.52 & -0.43 & 0.91 \\ 0.92 & -2.58 & 1.82 \\ -1.21 & 1.62 & -2.75 \\ 1.14 & -1.83 & 4.19 \end{bmatrix} \cdot h_t$$

$$y = W_{hy} \cdot \begin{bmatrix} 0.3 & 1.0 & 0.1 & -0.3 \\ -0.1 & 0.3 & -0.5 & 0.9 \\ 0.9 & 0.1 & -0.3 & 0.7 \end{bmatrix}$$

$$y \approx \begin{bmatrix} 1.0 & 0.5 & 0.1 & 0.2 \\ 2.2 & 0.3 & 0.5 & -1.5 \\ 3.0 & -1.0 & 1.9 & 0.1 \\ 4.1 & 1.2 & -1.1 & 2.2 \end{bmatrix}$$

Target chars:

"e" "i" "i" "o"

Date:

output layer

1.0	0.5	0.1	0.2
2.1	0.3	0.5	-1.5
-3.0	-1.0	1.9	-0.1
4.1	1.2	-1.1	2.2

All weight matrices

W_{hy} , W_{hh} , W_{wh}

$\uparrow \quad ? \quad \uparrow \quad \uparrow W_{hy}$ are randomly initialized

hidden layer

0.3	1.0	0.1	-0.3
-0.1	0.3	0.5	0.9
0.9	0.1	-0.3	0.7

? ? $\uparrow P_{wh}$

input layer

1	0	0	0
0	1	0	0
0	0	1	0
0	0	0	0

input chars "h" "e" "i" "i"

AFTER one-hot

• Recap of where we are, we just finish calculating our logits (output).

• In the very first time step we fed in h, where the RNN thinks the next character "h" is 1.0 likely

Logits: $\begin{bmatrix} 1.0 \\ 2.1 \\ -3.0 \\ 4.1 \end{bmatrix} \rightarrow \begin{bmatrix} "h" \\ "e" \\ "i" \\ "o" \end{bmatrix}$ To come next, "e" is 2.1
Step 1 likely to come next... and so on.

• RNN thinks "o" should come next! So it guessed wrong!

Date:

- These logits are raw scores, we need to apply Softmax to compute the logits into probabilities.

$$\text{Softmax}(Y_t)_i = \frac{e^{y_{t,i}}}{\sum_j e^{y_{t,j}}}$$

Where: $y_{t,i}$

- e is the exponent of the logit value

- $\sum_j e^{y_{t,j}}$ is the sum of all exponentiated logits in the column (normalization factor)

- Given logits.

$$Y = \begin{bmatrix} 1.0 & 0.5 & 0.1 & 0.2 \\ 2.2 & 0.3 & 0.5 & -1.5 \\ -3.0 & -1.0 & 1.9 & -0.1 \\ 4.1 & 1.2 & -1.1 & 2.2 \end{bmatrix}$$

- Compute softmax for each time step

$$\text{- Softmax for } T_1 = [1.0, 2.2, -3.0, 4.1]$$

1) Exponentiate each logit

$$e^{1.0} = 2.718, e^{2.2} = 9.025, e^{-3.0} = 0.050, e^{4.1} = 60.34$$

2) Sum exponentials $2.718 + 9.025 + 0.050 + 60.34 = 72.13$

3) Softmax = $\frac{[2.718, 9.025, 0.050, 60.34]}{72.13} = [0.037, 0.12, 0.007, 0.83]$
Probabilities

Softmax for $T=2$

- Logits = $[0.5, 0.3, -1.0, 1.2]$

1. Exponentiate each logit:

$$e^{0.5} = 1.649, e^{0.3} = 1.399, e^{-1.0} = 0.368, e^{1.2} = 3.320$$

2. sum of exponentials:

$$1.649 + 1.399 + 0.368 + 3.320 = 6.686$$

3. Softmax probabilities:

$$P_2 = \begin{bmatrix} 1.649 & 1.399 & 0.368 & 3.320 \\ 6.686 & 6.686 & 6.686 & 6.686 \end{bmatrix} = \begin{bmatrix} 0.2466 \\ 0.2019 \\ 0.0550 \\ 0.4965 \end{bmatrix}$$

• Repeat for $T=3$ and $T=4$

• Find softmax probability matrix:

$$P = \begin{bmatrix} 0.0377 & 0.2466 & 0.1131 & 0.1074 \\ 0.1251 & 0.2019 & 0.1687 & 0.0196 \\ 0.0007 & 0.0550 & 0.6641 & 0.0796 \\ 0.8365 & 0.4965 & 0.0341 & 0.7935 \end{bmatrix}$$

• each column sum to 1, verifying the softmax transform

Date:

- Compute the gradients of Loss w.r.t logs

$$\frac{\partial L}{\partial Y_t} = P_t - \text{target}_t$$

Where :

- P_t is The softmax probability matrix.

- target_t is the one-hot encoding of character

Given Targets

we are Training on The sequence "hell" $\rightarrow T$

t	Target character	one-hot-vector
1	"e"	[0, 1, 0, 0]
2	"l"	[0, 0, 1, 0]
3	"l"	[0, 0, 1, 0]
4	"o"	[0, 0, 0, 1]

- Compute $\frac{\partial L}{\partial Y_t}$ or gradients with Target Matrix T

$$T = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 \end{bmatrix}$$

- To compute $\frac{\partial L}{\partial y}$ we subtract the Target matrix from the softmax probability Matrix

$$\frac{\partial L}{\partial y} = P - \text{Target}$$

$t_1 \quad t_2 \quad t_3 \quad t_4$

$$= \begin{bmatrix} 0.0377 & 0.2466 & 0.1131 & 0.1074 \\ 0.1251 & 0.2019 & 0.1687 & 0.0196 \\ 0.0007 & 0.0550 & 0.6841 & 0.0796 \\ 0.8365 & 0.4965 & 0.0341 & 0.7935 \end{bmatrix} - \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 \end{bmatrix} \begin{matrix} t_1 \\ t_2 \\ t_3 \\ t_4 \end{matrix}$$

$$\frac{\partial L}{\partial y} = \begin{bmatrix} 0.0377 & 0.2466 & 0.1131 & 0.1074 \\ -0.8799 & 0.2019 & 0.1687 & 0.0196 \\ 0.0007 & -0.9450 & -0.3159 & 0.0796 \\ 0.8365 & 0.4965 & 0.0341 & -0.2065 \end{bmatrix}$$

$t_1 \quad t_2 \quad t_3 \quad t_4$

- This is our output layer error.

- It is our 3×4 matrix of hidden states (rows = 3 dim-hidden columns = 4 dim-out)

$$H = \begin{bmatrix} 0.3 & 1.0 & 0.1 & -0.3 \\ -0.1 & 0.3 & -0.5 & 0.9 \\ 0.9 & 0.1 & -0.3 & 0.7 \end{bmatrix}$$

- Next is the gradient formula for the hidden-to-output weights $W_{hy} \in \mathbb{R}^{4 \times 3}$ for each step x.

Gradient Formula for hidden-to-output weights

$$\frac{\partial L}{\partial W_{hy}} = \sum_{t=1}^4 \left(\frac{\partial L}{\partial y_t} \cdot h_t^T \right) = \sum_{t=1}^4 (4 \times 1) \times (1 \times 3) = 4 \times 3$$

- By collecting all 4 columns of $\frac{\partial L}{\partial y}$ (and H) into a single matrix, we are doing one single multiplication.
- Transpose H since h is a 3×4 is transpose is 4×3 .

$$H^T = \begin{bmatrix} 0.3 & -0.1 & 0.9 \\ 1.0 & 0.3 & 0.1 \\ 0.1 & -0.5 & -0.3 \\ -0.3 & 0.9 & 0.1 \end{bmatrix}$$

Thus

$$\frac{\partial L}{\partial W_{hy}} = \left(\frac{\partial L}{\partial y} \right) \times H^T$$

• Mult, by row by row.

- we get a final 4×3 result.

Denoted by : $[\Delta_{(1,:)}, \Delta_{(2,:)}, \Delta_{(3,:)}, \Delta_{(4,:)})]$

- Taking row 1 of \underline{H} which is
 $\underline{2y}$

$$[0.0377, 0.2466, 0.1131, 0.1074]$$

(we are making a 4×4 vector and multiplying it by H^T)

$$\Delta_{(1,:)} = [0.0377, 0.2466, 0.1131, 0.1074] \cdot \begin{bmatrix} 0.3 & 0.1 & 0.9 \\ 1.0 & 0.3 & 0.1 \\ 0.1 & -0.5 & 0.3 \\ -0.3 & 0.9 & 0.1 \end{bmatrix}$$

- computing each of the 3 entries.

Column 1 (the "(1,1)" element);

$$= 0.0377 \cdot 0.3 + 0.2466 \cdot 1.0 + 0.1131 \cdot 0.1 + 0.1074 \cdot (-0.3)$$

$$= 0.01131 + 0.2466 + 0.01131 - 0.03222$$

$$\approx 0.2370.$$

• Column 2 ≈ 0.1103

• Column 3 ≈ 0.0998

- Now row 1 The gradient is

$$\Delta_{(1,:)} \approx [0.2370, 0.1103, 0.0998]$$

• Row 2

• Row 3 → Similar for the following rows,

• Row

Date:

Final Gradient $\frac{\partial L}{\partial W_{hy}} \approx$

$$\begin{bmatrix} 0.2370 & 0.1103 & 0.0998 \\ -0.0496 & 0.0614 & -0.0841 \\ -1.003 & -0.0840 & 0.0566 \\ 0.8128 & -0.1376 & 0.6477 \end{bmatrix}$$

- Next we calculate The gradient with respect to h_t

- Once we have $\frac{\partial L}{\partial Y_t}$, we get $\frac{\partial L}{\partial h_t} = W_{hy}^T \left(\frac{\partial L}{\partial Y_t} \right)$

- This is another straight forward matrix multiplication.

- Then, we would compute The gradients w.r.t W_{hh}

- From $\frac{\partial L}{\partial h_t}$ via backpropagation through the recurrent const.

(Since our W_{hh} matrix is 0 we omit)

$$\frac{\partial L}{\partial W_{hh}} = \sum_t \left[\left(\frac{\partial L}{\partial h_t} \odot \frac{d \tanh(z_t)}{dz_t} \right) h_t^T \right]$$

- Similarly, we compute The gradients w.r.t W_{xh}

- From $\frac{\partial L}{\partial h_t}$ via backprop To The inputs \rightarrow hidden weights

$$\frac{\partial L}{\partial W_{xh}} = \sum_t \left[\left(\frac{\partial L}{\partial h_t} \odot \frac{d \tanh(z_t)}{dz_t} \right) x_t^T \right]$$