

## References:

why RNN? - <https://www.youtube.com/watch?v=4KpRP-YUw6c>

RNN architecture and forward propagation: [https://www.youtube.com/watch?v=BjWqCcbusMM&list=RDCMUCCWi3hpnq\\_Pe03nGxuS7isg&start\\_radio=1&rv=BjWqCcbusMM&t=958](https://www.youtube.com/watch?v=BjWqCcbusMM&list=RDCMUCCWi3hpnq_Pe03nGxuS7isg&start_radio=1&rv=BjWqCcbusMM&t=958)

RNN backpropagation: [https://www.youtube.com/watch?v=OvCz1acvt-k&list=RDCMUCCWi3hpnq\\_Pe03nGxuS7isg&index=6](https://www.youtube.com/watch?v=OvCz1acvt-k&list=RDCMUCCWi3hpnq_Pe03nGxuS7isg&index=6)

why do we need RNN?

RNN → Recurrent Neural Network

↳ is type of sequential model to work on sequential data, i.e., when sequence of data matters.

for example → text, speech

let's say we have a text classification problem, so given a sentence we should classify it. And we know that length of sentences vary, but ANN (Artificial Neural Networks) takes a fixed length vector as input. You might solve this problem by 0 padding all sentences to make the length of all sentences same as the sentence with maximum length. But this is not practical because of

↳ unnecessary computation

↳ Prediction problems

→ Sequence of words is being ignored.

Data for RNN

↓

(timestamps, input-features)

example: movie review → sentiment

	Review	x	Sentiment
$x_1$	movie was good	$x_{11}$ $x_{12}$ $x_{13}$	1
$x_2$	movie was bad	$x_{21}$ $x_{22}$ $x_{23}$	0
$x_3$	movie was not good	$x_{31}$ $x_{32}$ $x_{33}$ $x_{34}$	0

vocab - 5 words  
movie      was  
[1 0 0 0 0] [0 1 0 0 0] ...

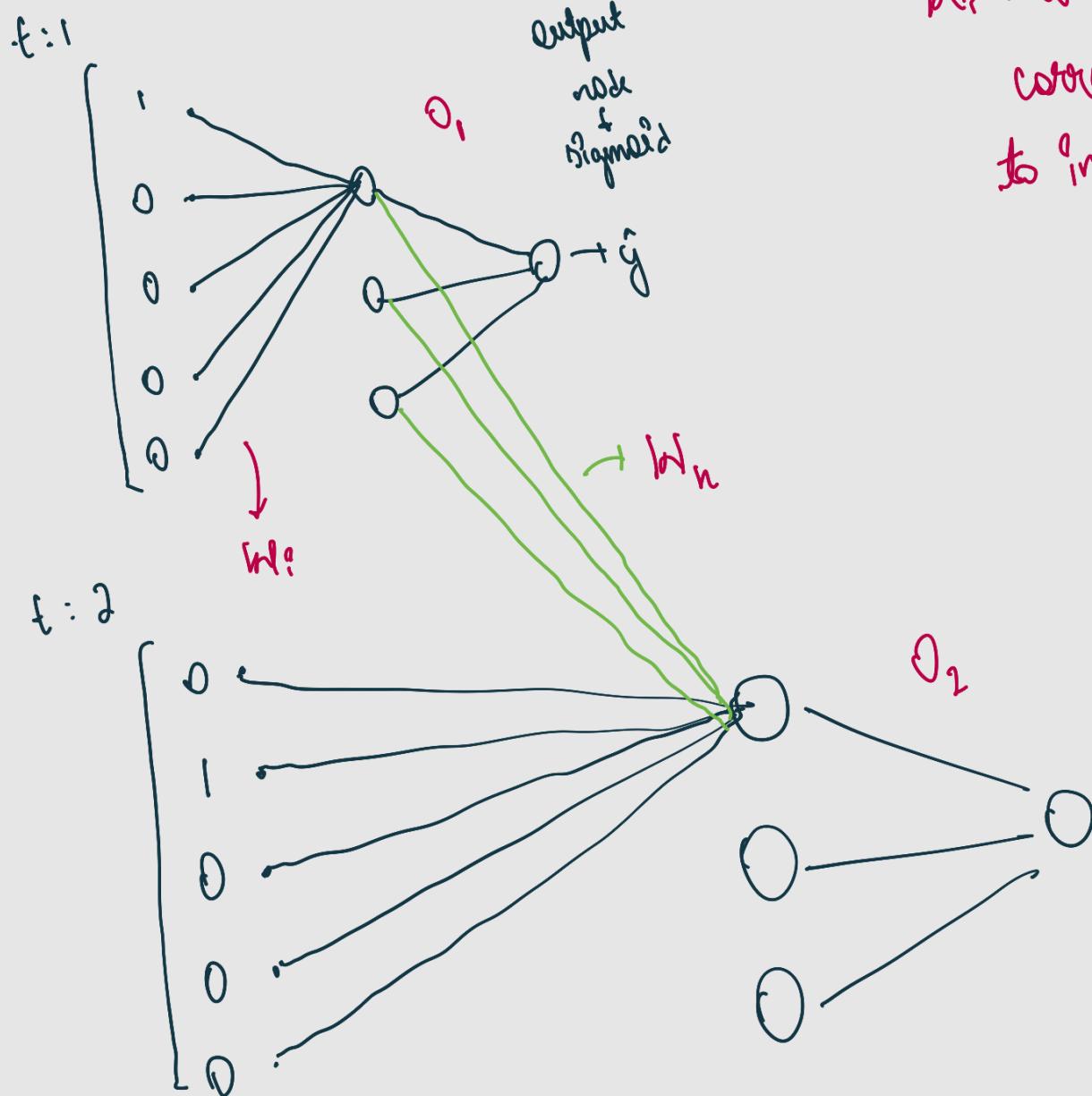
review<sub>i</sub> + movie was good  
↓

[{1 0 0 0 0} {0 1 0 0 0} {0 0 1 0 0}]  
↓  
(3, 5) # of features  
no. of time stamps

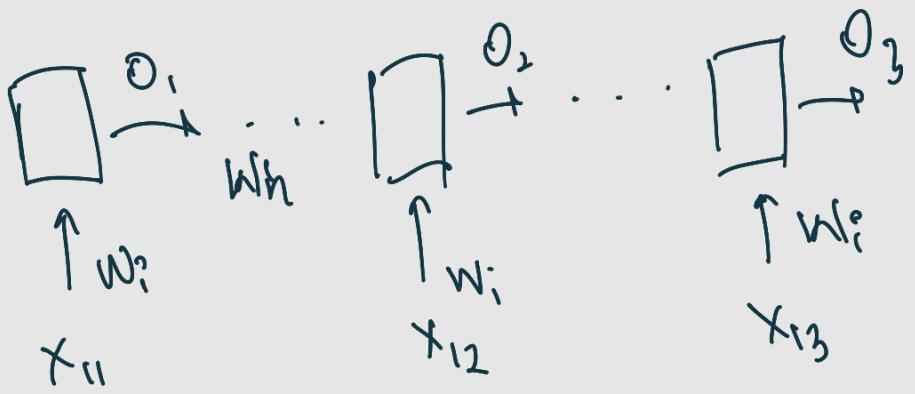
# Architecture



\* Data is not given as input in single timestamp

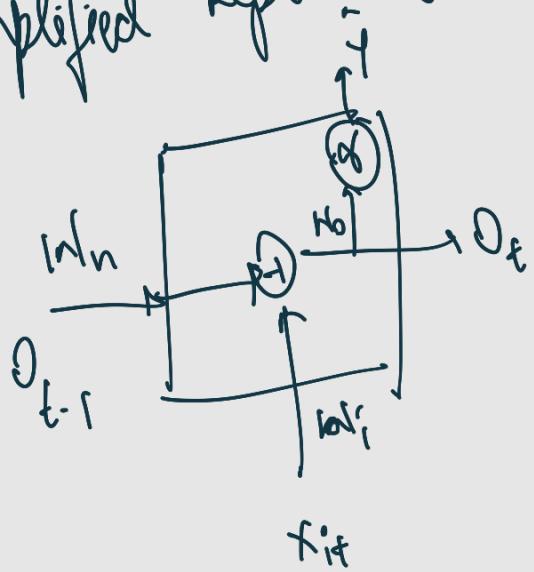


$w_i$  + weight  
corresponding  
to input feature



$$O_1: f(x_{i1}w_{i1}) \quad O_2: f((x_{i2}w_{i2} + O_1)w_{i3})$$

Simplified representation:



# Backpropagation in RNN

( $\uparrow$  single hidden)

let's take a toy dataset

Layer

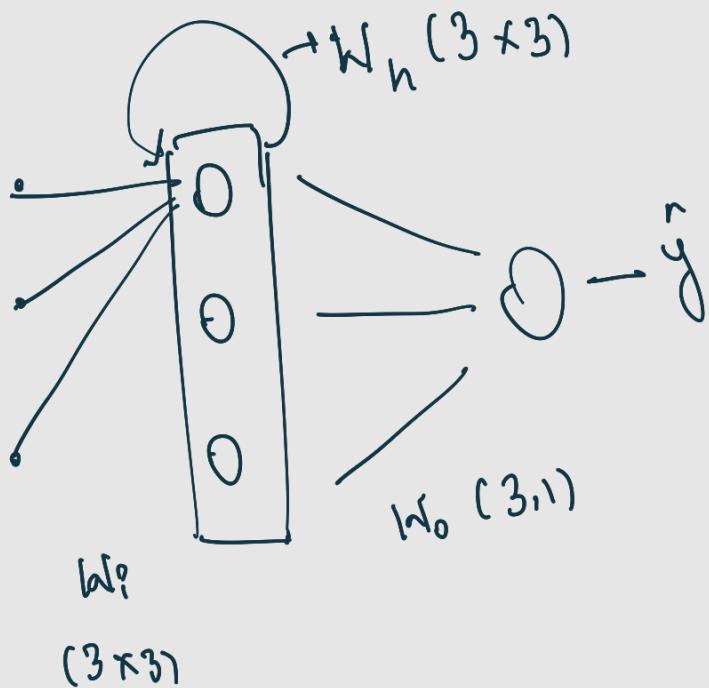
Many to one

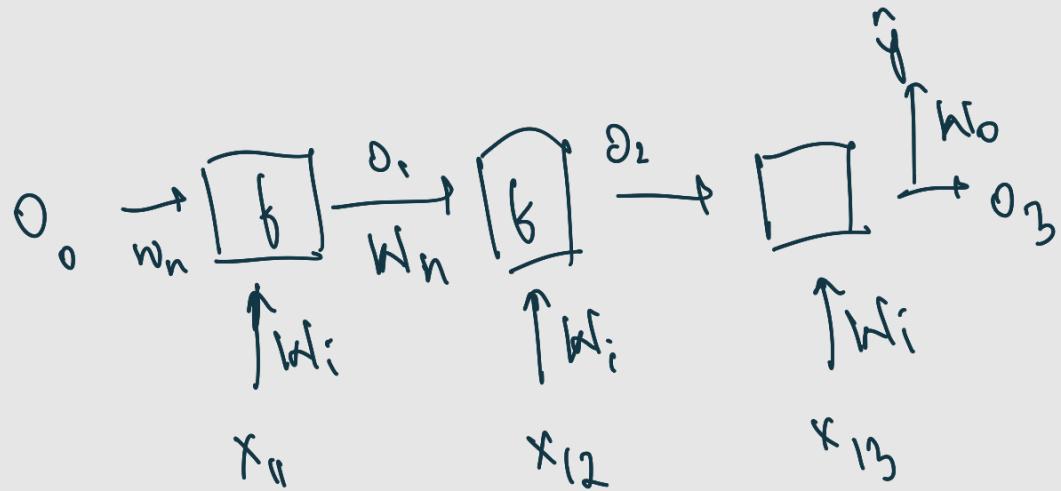
sentiment

text

cat	mat	rat	1
rat	mat	mat	1
mat	mat	cat	0

words  $\rightarrow$  cat mat rat  
 $[1 \ 0 \ 0] [0 \ 1 \ 0] [0 \ 0 \ 1]$





$$O_1 = f(O_0 w_n + b_0)$$

$$O_2 = f(O_1 w_i + b_1)$$

$$O_3 = f(O_2 w_h + b_2)$$

$$y - \hat{y} = \sigma(O_3 w_0)$$

$$L = -y_i \log \hat{y}_i$$

$$- (1-y_i) \log(1-\hat{y}_i)$$

learnable parameters

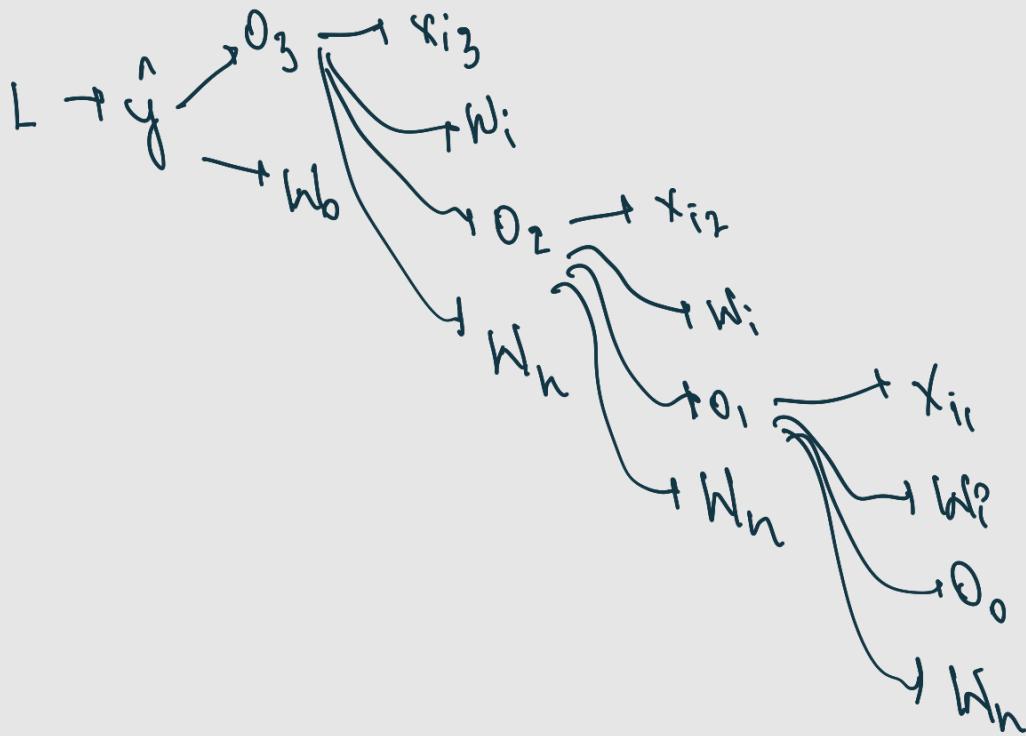
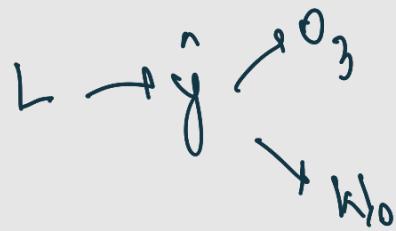
$$w_i \quad w_n \quad w_h$$

$$w_i^0 = w_i - \eta \frac{\partial L}{\partial w_i}$$

$$w_n^0 = w_n - \eta \frac{\partial L}{\partial w_n}$$

$$w_h^0 = w_h - \eta \frac{\partial L}{\partial w_h}$$

$$\frac{\partial L}{\partial w_0} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial w_0}$$



$$\frac{\partial L}{\partial k_{w_0}} : \frac{\partial L}{\partial y} \frac{\partial y}{\partial O_3} \frac{\partial O_3}{\partial w_0} +$$

$$\frac{\partial L}{\partial y} \frac{\partial y}{\partial O_3} \frac{\partial O_3}{\partial O_2} \frac{\partial O_2}{\partial w_0} +$$

$$\frac{\partial L}{\partial y} \frac{\partial y}{\partial O_3} \frac{\partial O_3}{\partial O_2} \frac{\partial O_2}{\partial O_1} \frac{\partial O_1}{\partial w_0}$$

$$\frac{\partial L}{\partial w_i} : \sum \frac{\partial L}{\partial y} \frac{\partial y}{\partial O_j} \frac{\partial O_j}{\partial w_i}$$

Aug for

$$\frac{\partial L}{\partial w_h} : \sum \frac{\partial L}{\partial y} \frac{\partial y}{\partial O_j} \frac{\partial O_j}{\partial w_h}$$