

Recurrent Neural Networks

A type of sequential model specifically designed to work on sequential data.

When we work with multiple samples whose order matters, we call that a sequence.

* ANN work well with tabular data.

* CNN work well with Images.

Reasons to not use ANN for Sequence Data:

1) Sequence data is of any length (with fixed input)

- text input — varying size

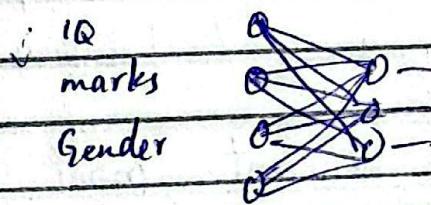
- If we apply zero padding (It would be unnecessary computation)

- Still if user input higher length of input it would be problem.

2) Sequence contains some meaning

- In ANN, all data input at once.

- It totally disregarding the sequence info.
IQ, marks, gender, Place



RNNs are a special class of neural network which have memory feature, which remember past inputs.

RNN comes in 1986 but famous recently.

RNN are usually thought of as a stepping stone to understanding fancier things like LSTM and Transformer.

Just like other neural networks, RNN have weights, biases, layers and activation functions.

- The big difference is that RNN also have a feedback loops.
- To deal with the problem of variable input size is to use RNN.

Types of Neural Networks (by structure)

1. Feedforward Neural Networks:

→ Information flows in one direction.
(input → output). No loops.

→ Feedforward Neural Networks just focus on current state.

$$x \rightarrow [H] \rightarrow y$$

Examples:

- MLP - Basic fully connected network
- CNN - Used for images & spatial data
- Autoencoders - For data compression and reconstruction,
- GANs (Generator) - Often built with FNN or CNN
- Transformers - Use feedforward layers with self-attention (not recurrent)

2. Recurrent Neural Networks:

→ It consider both the current input and past information by maintaining a hidden state that is updated overtime.

→ Previous outputs influence current processing. Designed for sequences.

Examples:

- Basic RNN
- LSTM - Handles long-term dependencies.
- GRU - Simpler alternative to LSTM.

Now to feed data to RNN (transforming text into numbers)
We input data to RNN in the below format.
(time steps, input-features)

Reviews	Sentiments
movie was good	1
movie was bad	0
movie was not good	0

The first step would be to convert text into vectors/numbers.

Integer Encoding

1) Form vocabulary:

↳ Find how many unique words document have.

→ Inputs are in English, our ML/DL algorithms did not understand it, we need to convert it into numbers/vectors.

↳ Give an integer value to every unique word.

→ Simple way is Integer Encoding.

↳ Replaces sentences with its integer values.

Padding

↳ Bring all to same size.

→ Size of each review is different.

vocab → 5 words

[10000] [01000] [00100] [00010] [00001]

Review 1: (3, 5)

Review 2: (3, 5)

Review 3: (4, 5)

Embeddings

→ Convert every word into a vector, on the basis of the whole document.

→ The vector here we got is dense not sparse.

→ Dense Values are non-zeros.

→ Sparse: Most of values are 0s.

Benefits:

• Dense representation

• Semantic Meaning (word use in which context)

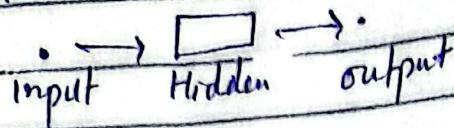
Word2vec and glove are word embedding techniques

thus is nice
↓
[0.7 0.1 0.3]

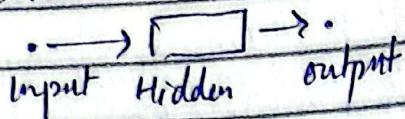
RNN Architecture:

→ RNN is very similar to ANN

ANN:



RNN:



→ But there two major differences.

1st:

How you fed data in your network?

→ In RNN input is sent on the basis of time.

time basis

$t=1$

$x_{11} \rightarrow [] \rightarrow$

$t=2$

$x_{12} \rightarrow [] \rightarrow$

$t=3$

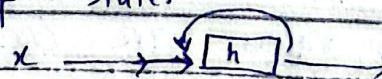
$x_{13} \rightarrow [] \rightarrow$

→ So first major difference is that complete input is not fed up at once in RNN.

2nd:

→ ANNs are feed forward neural networks. Information moves from input to output.

→ RNNs are not feed forward. In RNN there is a concept of state.



→ RNNs make use of an idea called state. This is just a description of a system (such as neural network) at any given time.

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

→ We represent every word as vector.

$[0000] [0100] [0010]$
 $[00010] [00001]$

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

...

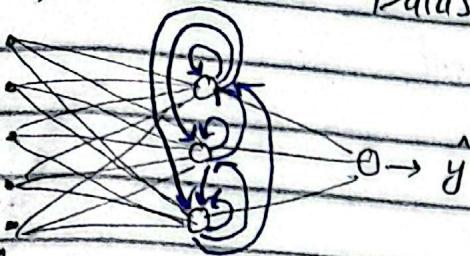
...

...

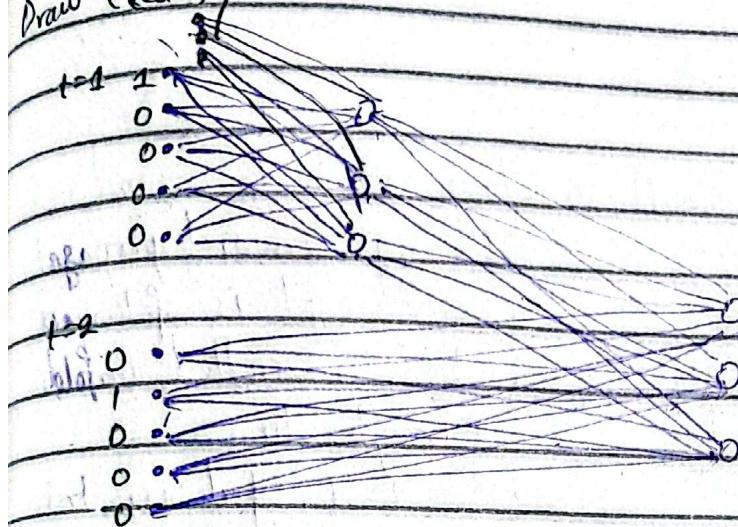
...

...

Architecture for Current Dataset:



Draw Clearly:



→ We are working with $t=1$ [10000]

→ It will go to hidden layer, some calculation done, activation will applied. Every node has an output, which will be sending back to the nodes, but for $t=2$.

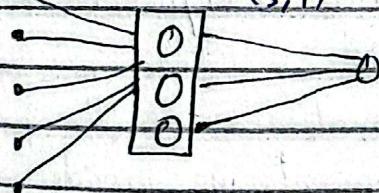
→ For $t=2$ [01000], sent to the same node but this time, node output come as well.

It will continue.

Number of Weights & Biases:

$$w = (5, 3)$$

$$(3, 1)$$



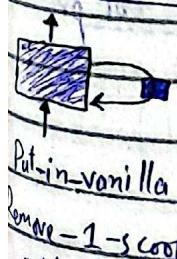
$$w_1 = 15 \quad w_2 = 3 \times 3 \quad w_3 = 3$$

$$\text{weights} = 27$$

$$\text{Biases} = 4$$

→ Now you may be thinking that for $t=2$, we got feedback from first, but for $t=1$ what will we get.

→ For $t=1$, feedback will be 0 or random values.



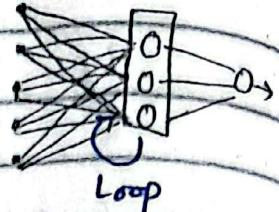
→ The loop on the right represent the state between one input and the next.

→ After each input, the system (represented by the big, light blue box) creates a new state, which goes into the black square.

→ This square is called the delay and we can think of it as a little piece of memory.

RNN Forward Propagation:

Review	Sentiment
$x_{11} \ x_{12} \ x_{13}$	1
$x_{21} \ x_{22} \ x_{23}$	0
$x_{31} \ x_{32} \ x_{33}$	0

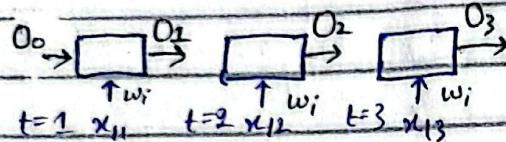


$x_{11} \rightarrow$ vector

$$[1 \ 0 \ 0 \ 0 \ 0]$$

↳ One by one we will sending words to neural network.

→ During forward propagation a concept is followed which is called unfolding through time.



$$1) f(x_{11}w_1 + O_0w_h + b_1)$$

→ This is done inside the above box.

→ On $t=1$, fed x_{11} take dot product with weight w_h → Apply activation function.

→ O_1 will be output from this node.

→ At $t=2$, we again use same network with the same weights, we'll feed 2nd word.

$$O_1 \rightarrow (1, 3)$$

$$w_h \rightarrow (3, 3)$$

$$2) f(x_{12}w_1 + O_1w_h + b_2)$$

$$(1, 5)(5, 3) + (1, 3)(3, 3)$$

$$(1, 3) + (1, 3)$$

$$O_2 \rightarrow (1, 3)$$

→ This time one thing is different, we provide an extra input to network which is O_1 .

$$3) f(x_{13}w_1 + O_2w_h + b_3)$$

$$O_3 \rightarrow (1, 3)$$

→ After $t=3$, our all words finished, take dot product of O_3 with three weights we get $(1, 1)$ which is scalar.

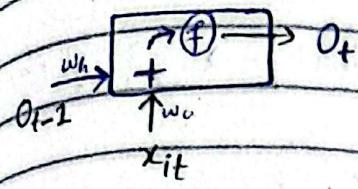
This is how forward propagation in RNN works.

(To keep consistency we fed

O_0 which is basically $(1, 3)$ vector of 0's or random)

→ Apply Sigmoid on scalar, you will get output.

Summary



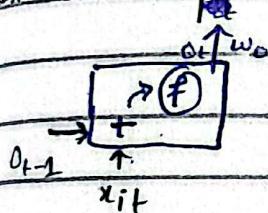
i = row number

t = time step

$$x_{it} \rightarrow$$

$$f(x_{it}w_i + O_{t-1}w_h)$$

Last Time Step:



$$\hat{y} = g(O_t w_o)$$

Types of RNN

Many to One One to Many Many to Many One to One

Many to One :

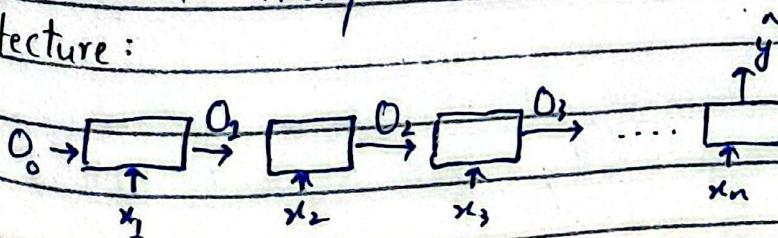
input : Sequence (sentence, time series)

output : non sequence (numbers → scalar (0, 1))

e.g.:

Sentiment Analysis

Architecture :



→ Imagine RNN as a box, where two input goes.

x_{it} and O_{t-1} (output of previous time step)

→ x_{it} goes through a weight matrix called w_i and O_{t-1} goes through a weight matrix w_h .

→ First perform dot products of input and weights, previous output and weights, then sum and then pass the vector to an activation function

You'll get current timestep output O_t .

→ If this is last timestep then pass it with w_o .

2) One to Many

input : non sequence
output : sequence

e.g: Image Captioning
Music Generation

Architecture:

```

graph LR
    x[x] --> NN[NN]
    NN --> y1[y1]
    y1 --> y2[y2]
    y2 --> ...[...]
    ... --> yn[yn]
    
```

Image → NN → A man is playing cricket.

3) Many to Many

input : Sequential Data

output : Sequential Data

- Same Length:

number of input = number of output

e.g: POS Tagging:

same is Imran

| It is also called
| Seq2Seq
| Same Length Variable
| length

e.g: POS Tagging:	My	name	is	Imran
	↓	↓	↓	↓
	pronoun	noun	verb	Noun

Architecture:

```

graph LR
    subgraph Input_Layer [Input Layer]
        x1[x1] --> y1[y1]
        x2[x2] --> y1
        ...
        xn[xn] --> y1
    end
    subgraph Hidden_Layer [Hidden Layer]
        y1 --> y2[y2]
        y2 --> y3[y3]
        ...
        ym[ym] --> ym
    end

```

- Variable Length:

e.g.: Machine translation

* Such architecture also called encoder-decoder.

* The first part, where you provide input to the encoder.

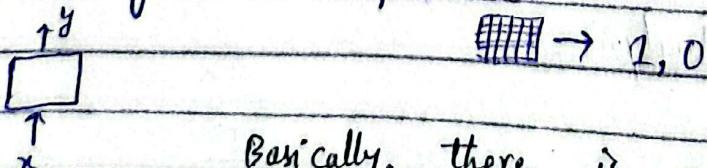
- * The part where you get output called decoder.

4) One to One

input : Non-Sequential data

output: Non sequential data

e.g.: Image Classification



Basically, there is no recurrence.
Technically, this is not RNN.

7 July 2025

Backpropagation in RNN

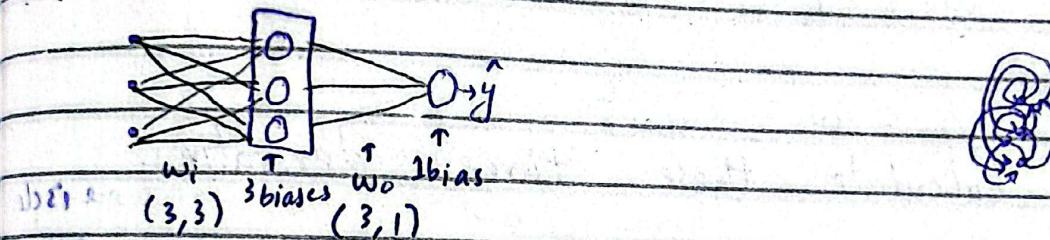
RNN is a neural network. Neural networks do learning. Learning is done through backpropagation. Here it is called BPTT.

dataset

	text			label	x	y
	cat	mat	rat	1	[1 0 0]	1
	rat	rat	rat	1	[0 0 1]	1
	mat	mat	cat	0	[0 1 0]	0

→ First convert text into vectors.

Vocab : eat mat rat
 $[1 0 0]$ $[0 1 0]$ $[0 0 1]$



Forward Propagation :

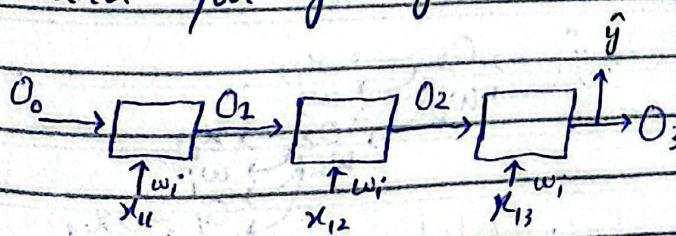
Reviews will be inserted word by word.

First send $[1 0 0]$, after calculations you get O_1 .

Then, send $[0 1 0]$, after processing you get O_2 .

Then, input $[0 0 1]$, after processing you'll get O_3 .

Since no words remaining in input, you move forward and you get \hat{y} .



$$O_1 = f(x_{11}w_i + O_0w_h)$$

$$O_2 = f(x_{12}w_i + O_1w_h)$$

$$O_3 = f(x_{13}w_i + O_2w_h)$$

When we get O_3 , we'll

use it to get \hat{y} .

$$\hat{y} = g(O_3w_o)$$

→ Now we got \hat{y} , we also get y . Use them to calculate loss.

$$L = -y \log \hat{y} - (1-y) \log(1-\hat{y})$$

→ After calculating loss, our goal is to minimize loss. The way to minimize loss is gradient descent algorithm. (optimization algorithm)
→ what we have to do is to find the values for parameters which reduce loss.

$$w_i = w_i - \alpha \frac{\partial L}{\partial w_i}$$

$$w_h = w_h - \alpha \frac{\partial L}{\partial w_h}$$

$$w_0 = w_0 - \alpha \frac{\partial L}{\partial w_0}$$

→ Now, calculate these three derivatives.

1. $\frac{\partial L}{\partial w_0}$

$$L \rightarrow \hat{y} \xrightarrow{O_3} w_0$$

→ This one is close to \hat{y} , so simple to calculate

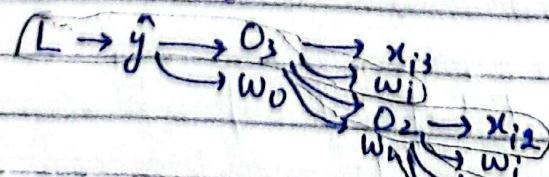
→ Change in loss depends on \hat{y} , because y is constant.

→ \hat{y} is calculated as $\hat{y} = \sigma(O_3 w_0)$

→ We can write using chain rule

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

2. $\frac{\partial L}{\partial w_i}$



→ L and w_i have a relationship of three separate paths.

There are three paths, so we have to calculate differentiation accordingly.

$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_3} \frac{\partial O_3}{\partial w_i} + \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_3} \frac{\partial O_3}{\partial O_2} \frac{\partial O_2}{\partial w_i} + \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_3} \frac{\partial O_3}{\partial O_2} \frac{\partial O_2}{\partial O_1} \frac{\partial O_1}{\partial w_i}$$

→ Here we have three terms because our RNN unfold in time, because our dataset have 3 words.

Compact form:

$$\frac{\partial L}{\partial w_i} = \sum_{j=1}^{t=3} \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_j} \frac{\partial O_j}{\partial w_i}$$

$$\Rightarrow \frac{\partial L}{\partial w_i} = \sum_{j=1}^n \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_j} \frac{\partial O_j}{\partial w_i} \quad n = \text{timesteps}$$

$$\frac{\partial L}{\partial w_h}$$

$$L \rightarrow \hat{y} \rightarrow w_0$$

$$\downarrow O_3 \rightarrow x_{i_3}$$

$$\downarrow w_i$$

$$\downarrow O_2 \rightarrow x_{i_2}$$

$$\downarrow w_h$$

$$\downarrow O_1 \rightarrow x_{i_1}$$

$$\downarrow w_i$$

$$\downarrow O_0$$

$$\downarrow w_h$$

→ L and w_h has a relationship of three separate paths.

$$\frac{\partial L}{\partial w_h} = \frac{\partial L}{\partial \hat{y}} + \frac{\partial \hat{y}}{\partial O_3} \frac{\partial O_3}{\partial w_h} +$$

$$\frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_3} \frac{\partial O_3}{\partial O_2} \frac{\partial O_2}{\partial w_h} +$$

$$\frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_3} \frac{\partial O_3}{\partial O_2} \frac{\partial O_2}{\partial O_1} \frac{\partial O_1}{\partial w_h}$$

Compact form:

$$\frac{\partial L}{\partial w_h} = \sum_{j=1}^n \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial O_j} \frac{\partial O_j}{\partial w_h} \quad n = \text{timesteps}$$

→ As these three derivatives calculated, apply gradient descent step.

→ First it work with some initial weights for w_i, w_h, w_0

Problems with RNN:

RNN → Sequential Data (text, time series)

↳ But RNN not used too much.

↳ It suffers from two major problems.

→ The more we unroll a recurrent neural network, the harder it is to train.

1. Problem of long-term dependency.

→ This problem arise because of vanishing gradient problem.

→ When a gradient becomes very small, learning slows down, and if a gradient becomes zero, learning stops entirely.

→ The more the timesteps are, the larger the term will be.

So your long-term dependency will get small.

The responsibility of derivatives comes on short term.

Solutions:

1) Use different activation functions

2) Better weight initialization.

3) Skip RNNs

4) LSTM

2. Problem 2: Unstable training (Exploding Gradient)

→ Gradients gets larger every time we step backward through the unrolled diagram.

Solutions:

1) Gradient Clipping

2) Controlled learning rate

3) LSTM

36