# Sequence Models

Examples of Sequence Data

1) Speech Recognition

2) music generation

3) Sentiment classification

4) DNA Sequence Analysis

5) Machine Translation

6) Video Activity Recognition

7) Named entity Recognition

Can be addressed as Supervised learning

$$X \longrightarrow Y$$

i/p     o/p

X and Y, can have Same or diff length.

# Notation

$x: \underline{\underset{x^{<1>}}{\text{Harry}} \ \underset{x^{<2>}}{\text{Potter}}} \ \text{and} \ \underline{\underset{x^{<3>}}{\text{Hermione}} \ \text{Granger}} \ \text{invented}$

new spell. $x^{<9>}$

NER — People name, Company's Name, Time, Location, currency Names, Country Names.

$$y: \quad 1 \quad 1 \quad 0 \quad 1 \quad 1 \quad 0 \quad 0 \quad 0$$

$$T_x = 9 \quad (\text{length of i/p sequence, } x)$$

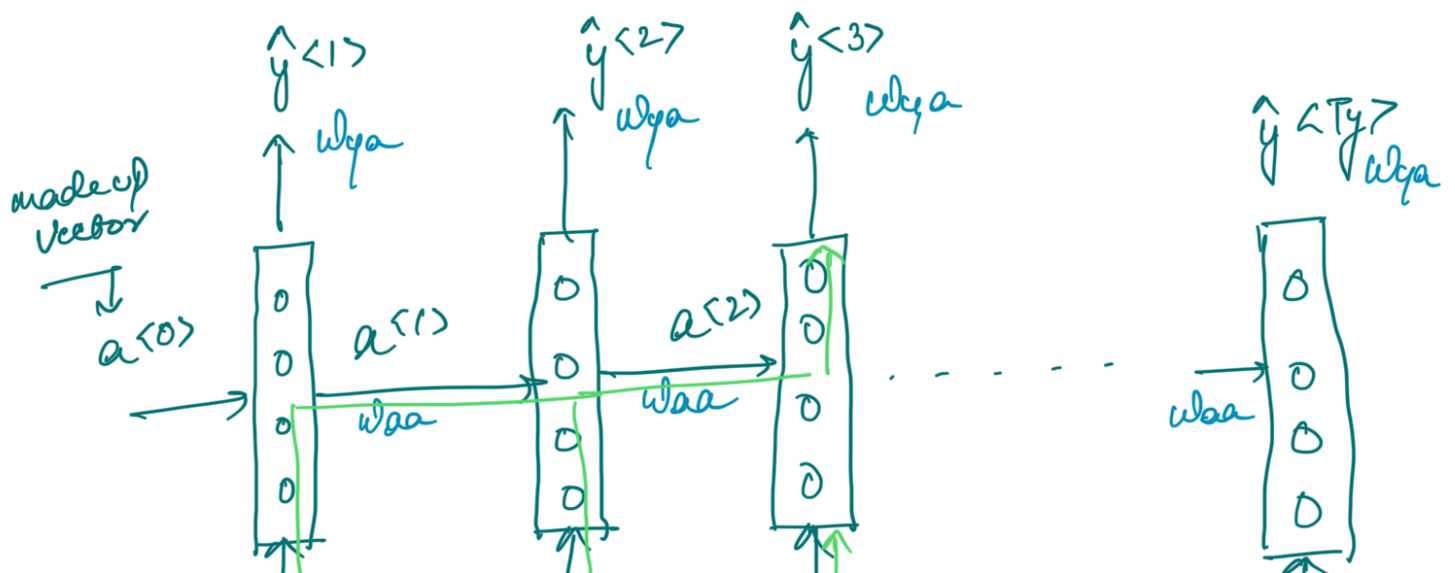→ Individual word representation ⟶ Vocabulary/
Dictionary

One-hot Vectors will be generated from the Sentence.

# RNN

why we cannot use Standard NN's:—

① I/p & o/p can be different lengths.

② Doesn't share features learned across different
position of text.

⇒ What is RNN

$X^{<1>}$ $W_{aa}$ (first word)

$X^{<2>}$ $W_{ax}$ (2nd word)

$X^{<3>}$ $W_{ax}$ . . . . .

$X^{<Tx>}$ $W_{ax}$ (4th word)

→ at every step activation will be passed on to the next layer.

→ RNN scans through the data from left to right.

→ Parameters at each time step are shared.

→ When predicting $\hat{y}^{<3>}$

$$x^{<3>} + x^{<2>} + x^{<1>}$$

→ **Weakness** : uses coords before it, not after.

$\hat{y}^{<3>}$ will only use till $x^{<3>}$ not $x^{<4>}$ or so on.

# Forward Propagation

usually **tanh** /ReLU

$a^{<0>} = \vec{0}$

$a^{<1>} = g_1(W_{aa} \, a^{<0>} + W_{ax} \, x^{<1>} + b_a)$

$\hat{y}^{<1>} = g_2(W_{ya} \, a^{<1>} + b_y)$ ⟵

$$0 \qquad 0 = \quad \text{o/p sigmoid activation}$$

$$\text{FP} \quad \Rightarrow \quad \begin{cases} a^{<t>} = g_1\left(\omega_{aa}\, a^{<t-1>} + \omega_{ax}\, X^{<t>} + b_a\right) \\[2mm] \hat{y}^{<t>} = g_2\left(\omega_{ya}\, a^{<t>} + b_y\right) \end{cases}$$

$$a^{<0>} \quad , \quad x^{<1>} \quad \Longrightarrow \quad a^{<1>} \rightarrow \hat{y}^{<1>} \quad \cdots$$

# Backward Propagation :-

$\hat{y}^{<1>}$    Wya

$\hat{y}^{<2>}$    Wya

$\hat{y}^{<3>}$    Wya

$\hat{y}^{<T_y>}$    Wya

made up vector

$a^{<0>}$

$a^{<1>}$   Waa

$a^{<2>}$   Waa

Waa

Waa

$X^{<1>}$   Waa

$X^{<2>}$   Wax

$X^{<3>}$   Wax

$X^{<T_x>}$   Wax

( first word )

( 2nd word ) . . . .

( 4th word )

$\longleftarrow$ : Backprop.

# You need loss function for back Prop.
  Standard logistic Regression loss

$$L^{<t>}(\hat{y}^{<t>}, y^{<t>}) = y^{<t>} \log \hat{y}^{<t>} - (1 - y^{<t>}) \log (1 - \hat{y}^{<t>})$$

for 1 time-step

$$L(\hat{y}, y) = \sum_{t=1}^{T_y} L^{<t>}(\hat{y}^{<t>}, y^{<t>}) \Rightarrow \text{loss of all}$$

# Types of RNNs :

$$T_x = T_y$$

i/p        o/p
        ↓

→ many-to-many
   (NER)

→ many-to-one
   (Sentiment classification)

→ one-to-many
   ( music generation)

→ many-to many ($T_x \neq T_y$) like machine Translation.

also called as <u>encoder – decoder</u>.

# <u>Language Model and Sequence Generation</u>.

→ Most basic and imp task in NLP → Language Modelling.

<u>What is Language Model :-</u>

Speech recognition example

The apple and pair salad.        ①

✓ The apple and pear salad        ②

                        using Lang Model

way ② is picked ⇒ by using ~ ~ ~ f ~

↓ feels what is prob of each Sentence.

$$① \longrightarrow 3.2 \times 10^{-13}$$

$$② \longrightarrow 5.7 \times 10^{-10} \checkmark$$

# A <u>Language Model's job</u>
→ what is the probability of the Sentence.

$$P(\text{Sentence}) = ?$$

$$P(y^{<1>}, y^{<2>}, \dots y^{<T_y>})$$

→ How to build <u>LM</u> using RNN :—

Training set : large corpers of english text.

<u>Step 1</u> :- Tokenize (Vocab)

<u>Step 2</u> :- One-hot Vectors / indices, Also,

add extra token <u>{EOS}</u> : end of Sentence.

<UNK> : for unknown token

eg Sentence:- Cats average 15 hours of sleep a day. <EOS>

$P(a)\ P(\text{aa})\ P(\text{aaron})\ldots P(\text{Cats})$
$$\hat{y}^{<1>}$$

$P(\text{average} \mid \text{Cats})\ \ldots$
$$\hat{y}^{<2>}$$

$P(<EOS> \mid \ldots\ldots)$

$y^{<1>}\ \ y^{<2>}\ \ y^{<3>}\ \ :\ \ 3\ \text{word Sentence.}$

$$= P(y^{<1>}) * P(y^{<2>} \mid y^{<1>}) * P(y^{<3>} \mid y^{<1>}, y^{<2>})$$

Probability of 3 word Sentences.