

* (Course 5) Sequence models

(week 1)

→ why sequence models

- models like RNNs have transformed speech recognition NLP & others

- in speech recognition we give input $X(MMM)$ → "Hi"

- Music generation is another example with seq data

The input may be (empty set, single int. referring to genre) output is sequence

- sentiment classification input phrase → classify it

- DNA seq analysis from given seq. label which part is protein

- Machine translation moi → Me

- Video activity recognition input (seq. of frames) → running

- Name entity recognition omar is --- → omar.

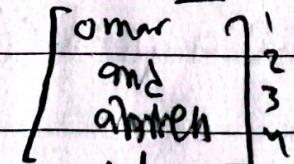
it can also be used find people's names, companies names times, locations, countries names, currency names

→ Notation

- $X: \overset{x^{(1)}}{o} \overset{x^{(2)}}{m} \overset{x^{(3)}}{a} \dots$ $T_x = \text{length}$

$y: \overset{y^{(1)}}{l} \overset{y^{(2)}}{o} \overset{y^{(3)}}{m} \dots$ $T_y = \text{length}$

- to represent a word in sentence first we come up with a vocabulary or dictionary.



then use one-hot representations to represent each word

→ Recurrent Neural Network model.

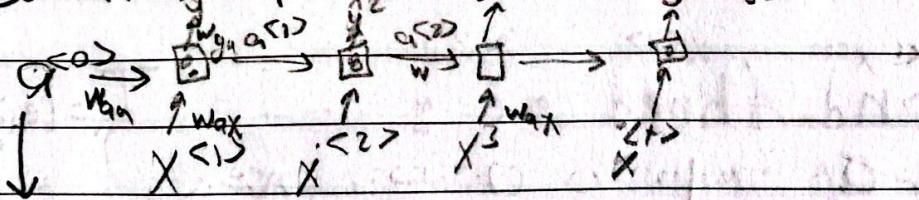
- we will talk about how you can build model, built NN to learn the mapping from $X \rightarrow y$ → Bec RNN handle seq. data & have
- why not standard NN. memory that capture dependencies between input & output
- ① Doesn't learn features across different position of text
- ② inputs, outputs can be different lengths

• in RNN if we are reading a sentence from left to right. The first word we will read is $x^{<1>}$ then feed it to NN then the NN try to predict if it is a person's name or not

$$x^{<1>} \rightarrow \boxed{0} \rightarrow y^{<1>}$$

• what RNN does is when it goes on the road the second word in the sentence instead of just predicting y_2 using x_2 only it also gets to input some information from that timestep one so in particular, deactivation value from time step one is passed on

to time step two $a^{<2>} \quad g^{<2>} \quad g^{<1>}$



vector of zeros

• one weakness of RNN is that it uses the info that is earlier in the seq. to make a prediction, but not use info later in the sequence

$$\begin{aligned} a^{<1>} &= g(w_{ha} a^{<0>} + w_{ax} x^{<1>} + b_a) \leftarrow \text{Tanh/ReLU} \\ y^{<1>} &= g(w_{ya} a^{<1>} + b_y) \leftarrow \text{sigmoid/Softmax} \end{aligned}$$

$$\begin{aligned} a^{<2>} &= g(w_{ha} a^{<1>} + w_{ax} x^{<2>} + b_a) \\ y^{<2>} &= g(w_{ya} a^{<2>} + b_y) \end{aligned}$$

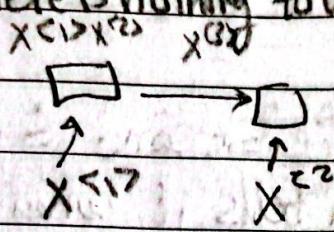
→ Backpropagation Through time

→ Different types of RNN

- Sentiment Classification: $x = \text{text}$

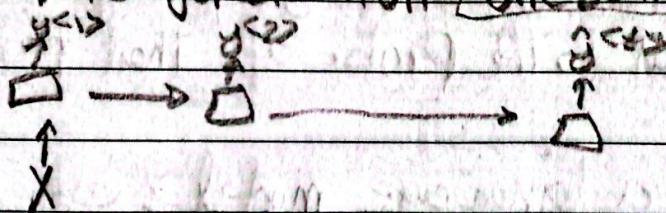
\xrightarrow{x} num from 1-5
or 0/1
 \downarrow
Sequence

Ex: There is nothing to like in movie +ve or -ve review



many to one

- Music generation [one to many]



- Machine translation many to many but input & output could be different length.

it has encoder which take as input the french sentence & then the decoder which having read in the sentence & outputs the translation into a different language

→ language model & sequence generation.

language modeling is one of the most basic & important tasks in NLP. we will talk about building a language model using an RNN.

- what is lang. model? lets say we are building a speech recognition system & you hear the sentence

→ The apple & ~~pear~~ Salad was delicious
(not pair)

→ the way SR sys. picks the second sentence is by using a language model which tells it what is the prob. of either of these two sentences.

• so what lang. model does is given any sentence its job is to tell you what is the prob. of that particular sentence

- How do you build a language model? using RNN
- first need a training set (~~corpus~~) of english text or whatever language

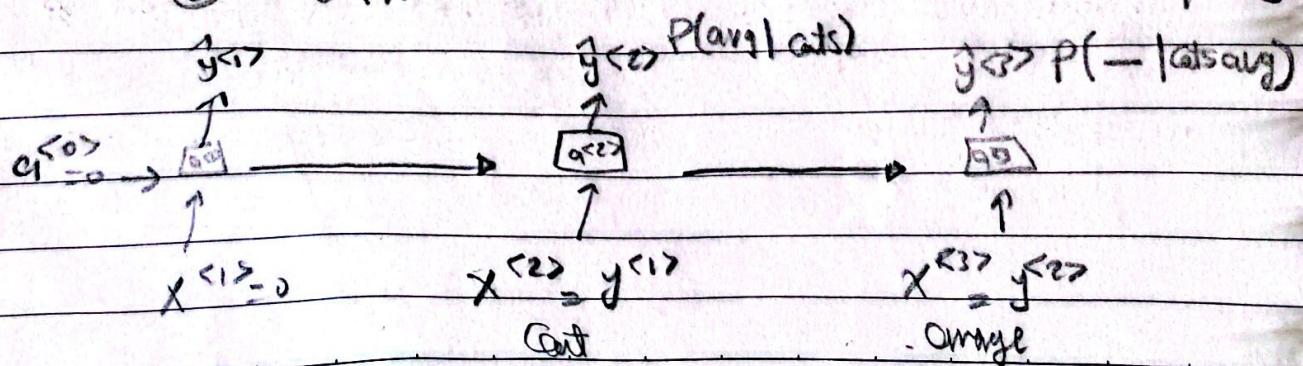
→ in nlp means tens of english sentences

- ex a sentence in your training set as follows:

Cats Average 15 hours of sleep a day

① tokenize the sentence (Form a vocab & map each word to one hot vector)

② build RNN to model the chance of these diff. sequences



- This RNN learns to predict one word at a time going from left to right.

→ Sampling Novel Sequences ✕

- After training a sequence model, one of the ways you can informally get a sense of what is learned is to have a sample novel sequences

- Using a character level language model has some pros & cons one is that you don't ever have to worry about unknown word tokens , in a character level language model is able to assign a sequence like man, ex. word if man is not in the vocabulary for the word level lang. model you just have to assign it the unknown word token , But the main disadvantage of the character level lang. model is that you end up with much longer seq.

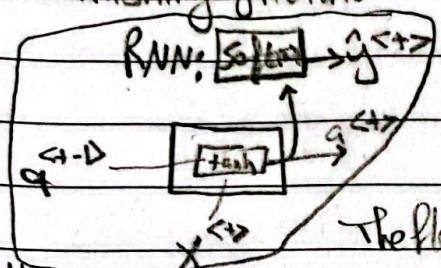
→ Vanishing Gradient with RNNs

- for ex: the cat, which already ate ----, was full
"Cats, --, --, were full"
- This is example that language can have long-term dependencies where earlier words can affect what needs to come later
- it might be difficult to get a NN to realize that it need to memorize a singular noun @ plural noun @ to that later in seq. it can generate either was @ were

if derivitivies do explore we apply Gradient clipping

→ Gated recurrent unit (GRU)

- It is modification to RNN hidden layer that makes it much better at capturing long-range connections & helps along with vanishing gradient.



- The GRU have C → memory cell that provide memory that stores info to remember it over long time

- F_t → Gate beta -1 usig sigmoid. It controls

The flow of info through the network. Update gate (decides how much of prev. memory to keep & how much to update with new info). Reset gate (decides how much to forget before updating the hidden state with new input)

→ Long short term memory (LSTM)

- More powerful than GRU and allows us to learn long term seq.

- LSTM is type of RNN designed to Solve Vanishing gradient and learning over long sequences

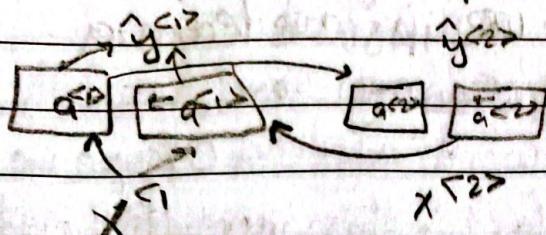
- LSTM introduces a memory cell with three gates

~~update~~ ← input gate : decides how much new info to store in memory cell
 forget gate: // " " " info "forget" // " "
 output gate: Controls the output based on the memory //

- memory cell is core component that stores the info. Gates are mechanism that controls the flow of info to/fout the memory cell & The gates used AF (sigmoid) to decide how much the info to allow through

→ Bidirectional RNN

- allows us at the point in time to take info from both earlier & later in sequence



- The disadvantage is that we need the entire seq. of data before you can make predictions anywhere.

In BiRNN, one layer processes the input from start to end (forward) & other layer processes the sequence from end to start (backward). At each time step, the outputs from the forward & backward are combined to form the final output.

- GRU, LSTM, SimpleRNN are types of recurrent NN used in language models but they are not considered seq2seq models.