# Notes

26 March 2024   07:40

## Lecture Outline:



## Simple RNN-based Language Model:  Recap



### Training an RNN-based LM:

- Dataset: A big corpus of words

- Model: RNN
    - Feed a set of words (sentence if you will) to the model
    - Model predicts probability distribution (of next word) over the corpus

- Loss function:
    - Cross entropy between GT prob of next word and Pred prob at each time step
    - Total loss in one iteration: Avg. Of loss at each time step



- Teacher Forcing:
    - Calculating loss at each time step allows model to predict only the next word.

- ○ Another way could be: Give model couple of words as input e.g. "the students" and let the model finish the sentence. The loss can be calculated over the sentence vs only one word.

- Backpropagation and Gradient for RNNs:
  - ○ (if it is not clear until now) There is only one/shared weight matrix $W_h$.
  - ○ However (downstream )gradient at each timestep t will be different as upstream (gradient) is different.
  - ○ What to do? Sum all the gradients at all timesteps and apply all of them at once (Do not chain the updates as during forward pass each $W_h$ was used as a separate matrix)



# Other Uses of RNNs:

## Text Generation:

- Notice how output at time t is being fed back as input of time t+1.



## Sequence Tagging:
- Assigning Part of speech to each word in the sentence



## Sentiment Analysis:
- Using RNN to get encoding of sentence and use a classifier layer to get a score

## Speech-to-Text/Text Generation:

- Using a NN as conditioning layer (to transform audio signal into embedding vectors)



# Evaluating RNN LM:

## Perplexity:
- Standard eval metric.
- Geometric representation of inverse probability.
- Indicates how well the model predicted the next word? In Ideal case, the model should have perplexity of 1.
- It is a unitless measure. Ranging from 1 to positive inf.
- e.g. perplexity value of 53 sort of means that model is choosing the next word from 53 near equiprobable words (out of all output prob. Values  53 values are very close to each other) -- representing model's uncertainty.



# Gradient Problems: Vanishing and Exploding

## Vanishing Gradients:

- No update at all to the model weights.

- Why the gradients vanish?
  - Visible reason – if the partial derivatives of hidden stages are small enough in magnitude, the upstream gradient becomes small and small as it gets propagated.
  - Architectural Reasons(not talked in the lecture):
    - Reuse of weights: during initialization/training if weights are small enough then upstream gradients will become smaller and smaller at each time step
    - Choice of Activation function: Commonly used tanh/sigmoid have saturation regions where the gradient is very small.
  - What is too small gradients/weights?
    - If eigen values of gradients/weights matrix are < 1 then it can be considered as too small

- Why is this a problem?
  - Model is rendered incapable of understanding long-term dependencies – as upstream gradient vanishes, timesteps in the beginning receive almost zero gradients and then the weights are not modified for the errors in those time steps.

## Exploding Gradient:
- Weight updates becoming too big

- Solution: Make gradients manageable (not too big) -- Gradient Clipping



**Gradient clipping: solution for exploding gradient**

- Gradient clipping: if the norm of the gradient is greater than some threshold, scale it down before applying SGD update

**Algorithm 1** Pseudo-code for norm clipping
$$\hat{g} \leftarrow \frac{\partial \mathcal{E}}{\partial \theta}$$
$$\text{if } \|\hat{g}\| \geq threshold \text{ then}$$
$$\hat{g} \leftarrow \frac{threshold}{\|\hat{g}\|} \hat{g}$$
$$\text{end if}$$

- Intuition: take a step in the same direction, but a smaller step

- In practice, remembering to clip gradients is important, but exploding gradients are an easy problem to solve

Source: "On the difficulty of training recurrent neural networks", Pascanu et al, 2013. http://proceedings.mlr.press/v28/pascanu13.pdf

## LSTM:

- As hidden states are almost completely overwritten at each time step --
  - RNNs not able to retain information for long
  - Also contributes to vanishing gradients – frequent changes might zero out some elements of hidden states

- How about adding some memory to RNNs so that hidden states can be changed selectively?



**4. Long Short-Term Memory RNNs (LSTMs)**

- A type of RNN proposed by Hochreiter and Schmidhuber in 1997 as a solution to the vanishing gradients problem.
  - Everyone cites that paper but really a crucial part of the modern LSTM is from Gers et al. (2000) 💜

- On step $t$, there is a hidden state $h^{(t)}$ and a cell state $c^{(t)}$
  - Both are vectors length $n$
  - The cell stores long-term information
  - The LSTM can read, erase, and write information from the cell
    - The cell becomes conceptually rather like RAM in a computer
- The selection of which information is erased/written/read is controlled by three corresponding gates
  - The gates are also vectors length $n$
  - On each timestep, each element of the gates can be open (1), closed (0), or somewhere in-between

"Long short-term memory", Hochreiter and Schmidhuber, 1997. https://www.bioinf.jku.at/publications/older/2604.pdf
"Learning to Forget: Continual Prediction with LSTM", Gers, Schmidhuber, and Cummins, 2000. https://dl.acm.org/doi/10.1162/089976600300015015
42

- Graphical representation of LSTM Block :



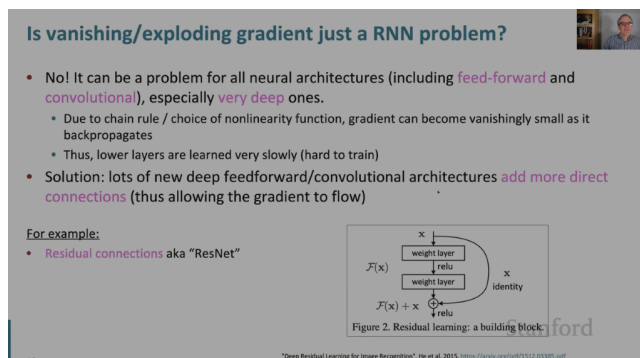**Long Short-Term Memory (LSTM)**

- Mathematical representation of LSTM Block



- Key difference to simple RNNs:
  - In LSTM: Not all is changed (added/forgotten) in each iteration
  - When new information is added to hidden states (in simple RNN) it is added by multiplication vs addition when new information is added to Cell states in (LSTMs)
  - In RNNs – multiplication carries out drastic changes vs in LSTMs addition carries out gradual changes
  - This makes retaining information lot easier for LSTMs

- How LSTMs help to control/reduce vanishing gradients problem?
  - By controlling rate of change of information in cell states, LSTM create a longer term version of short term memory
  - This prevents frequent changes to hidden states allowing upstream gradient to flow for until timestamp 0 for longer duration of time
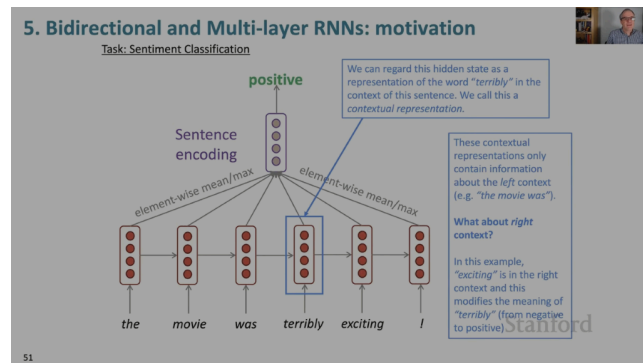
## Vanishing/Exploding Gradients in general DNN:

- Not only RNN specific. Every Deep Network suffers from such issues due to chaining the upstream gradient for longer depths
- Solution:



## Bi-directional and Multi-layer RNNs:

- Adding bi-directional context to RNNs: not only left context (from words on the left/previous words) but also right context (from words on the right of the current word)
- Only applicable when you have access to entire sentence e.g. when in sentiment analysis. Not applicable in e.g. language

modelling



- How?? -- Add a second backward RNN and concatenate both RNN's hidden states