

Input: Text, Task: Predicting next word in a sentence

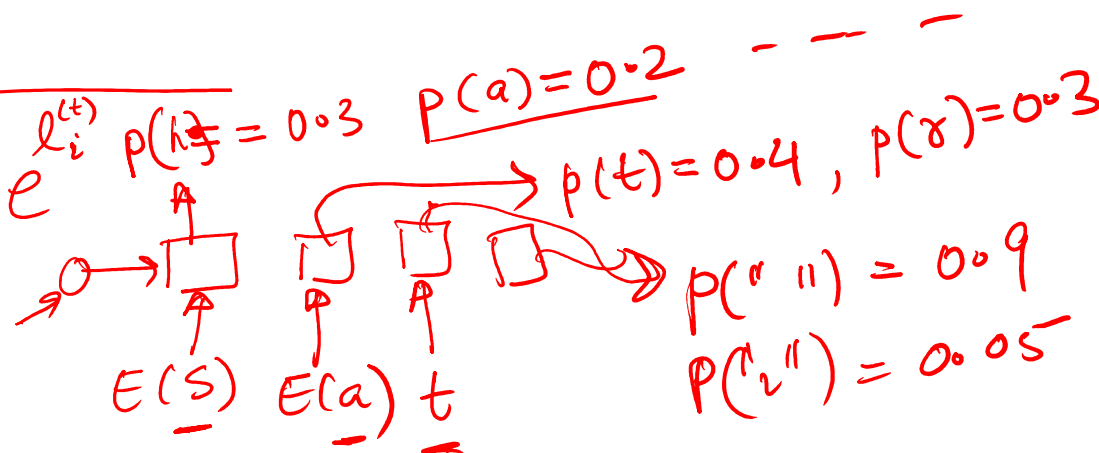
- Input sequence: x_1, x_2, \dots, x_n
- Each x_t instead of a real-valued vector is a discrete word.
- The task is to predict the following word $x_{\{n+1\}}$
- Output sequence: y_1, y_2, \dots, y_n can be written as $y_t = x_{t+1}$

$\boxed{h^{(1)}}$, $\underline{h^{(2)}}$ — $h^{(n)}$ — RNN states
 $P^{(1)}(\hat{y}_1)$ — \downarrow need to convert $\xrightarrow{(n)}$ $P(y_n)$ ← probability distribution over characters.

last linear layer parameters $\rightarrow h^{(t)} U = \underline{l^{(t)}} \in \mathbb{R}^V$
 for example: $\underline{l^{(t)}} = [-0.1, 2, -3, \dots]$ $|V|$

$$P^{(t)}(y_t = j | l^{(t)}) \equiv \frac{e^{l_j^{(t)}}}{\sum_{i=1}^{|V|} e^{l_i^{(t)}}}$$

$$-\log P^{(t)}(y_t = x_{t+1} | \underline{l^{(t)}})$$



The sequence prediction task

- Given a complex input \mathbf{x}
 - Example: sentence(s), image, audio wave
- Predict a sequence \mathbf{y} of discrete tokens y_1, y_2, \dots, y_n
 - Typically a sequence of words.
 - A token can be any term from a huge discrete vocabulary
 - Tokens are inter-dependent
 - Not n independent scalar classification task.



Translation

Context: \mathbf{x}

Predicted sequence: \mathbf{y}

Where can I find healthy and traditional Indian food? →

स्वस्थ और पारंपरिक भारतीय भोजन कहाँ मिल सकता है?

- Pre-DL translation systems were driven by transfer grammar rules painstakingly developed by linguists and elaborate phrase translation
- Whereas, modern neural translation systems are scored almost 60% better than these domain-specific systems.

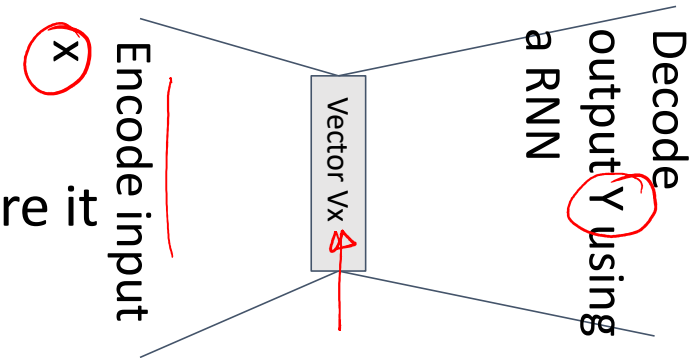
Sequence to Sequence learning



- Input sequence \mathbf{x}
 - Example: sentence(s), audio wave
- Predicted sequence \mathbf{y} has discrete tokens
 - Typically a sequence of words from a huge discrete vocabulary
- A common Model for Seq2Seq learning --- encoder-decoder model.

Encoder Decoder Model

- Encode x into a fixed-D real vector X
- Decode Y a token at a time by conditioning the generation of the next token on the tokens before it using a RNN

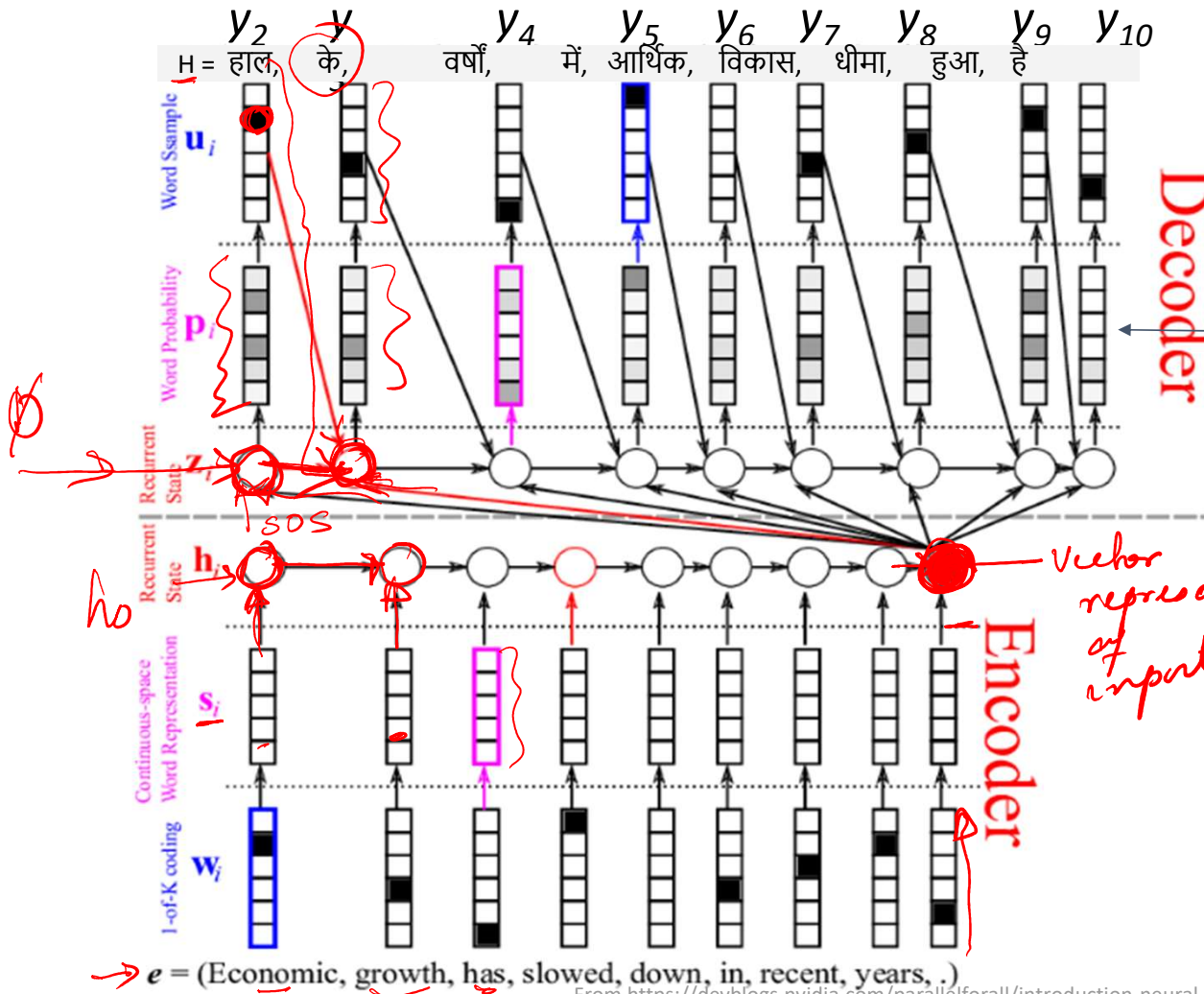


$$P(\mathbf{y}|\mathbf{x}, \theta) = \prod_{t=1}^n P(y_t | y_1, \dots, y_{t-1}, \mathbf{x}, \theta)$$

Handwritten notes:

- θ parameters
- y_1, y_2, \dots, y_T
- chain rule of probability
- Auto-regressive model

Encoder-decoder for sequence to sequence learning



Predicted sequence:

Choose high probability token and feed to next step.

$$\Pr(y_t | y_1, \dots, y_{t-1}, \mathbf{x})$$

RNN to generate y

RNN e.g. LSTMs to summarize x token-by-token

Embedding layer to convert each word to a fixed-D real vector

Context: \mathbf{x}

From <https://devblogs.nvidia.com/parallelforall/introduction-neural-machine-translation-gpus-part-2/>

Attention

- Attention
 - [Visualizing A Neural Machine Translation Model \(Mechanics of Seq2seq Models With Attention\) – Jay Alammam – Visualizing machine learning one concept at a time. \(jalammar.github.io\)](http://jalammar.github.io)

Inference

- Greedy algorithm
- Beam-search

Transformers

- RNNs require sequential computation of state, which makes training slow
 - Transformers replace “state” with multi-layered self-attention.
 - All positions in a sequence can be trained in parallel.
-
- <https://jalammar.github.io/illustrated-transformer/>