

## Dealing with Texts

- Words are categorical objects (non-continuous)

Suppose: ["the", "quick", "brown"]

$x_1 = \text{"the"}$

$x_2 = \text{"quick"}$

$x_t$  not numerical

$$h_t = \sigma(W_{xh}^T x_t + W_{hh}^T h_{t-1} + \beta_h)$$

## Problem with One-hot encoding:

- Need to create Array representing all english words  
↳ All values 0 except for 1.
- Input =  $T \times V$  ← size of array containing all Eng words (> 200 000 words)  
↑  
No. words
- All 1-hot encoded words have euclidean distance of  $\sqrt{2}$   
↳ due to '1'

## Solution: Embedding Layer

↙ opposite is sparse vectors (most elements are 0)

- represent words, tokens, textual elements as dense vectors
- vector is low-dimensional  
↳ instead of 10,000 for size 10,000 Vocab  
↳ size 50, 100, 300 instead.

### Step 1: Convert words to integers

↙ represent index for Embedding matrix

["I", "Like", "cats"] → [50, 25, 3] (between 0 & vocab-size) :  $E[50], E[25], E[3]$

### Step 2: Integers represent the index to retrieve corresponding dense vector from Embedding Matrix

Embedding Matrix  $E = (\text{Vocab-size} \times \text{embedding-dim})$

unique tokens      no. of dim  
eg: 10 000              eg: 50

$$E = \begin{bmatrix} 0.11, -0.34 \dots x_{50} \\ \vdots \\ 0.45, 0.67 \dots x_{50} \end{bmatrix}, \begin{matrix} \text{50 dim.} \\ \text{- row 0} \\ \text{- row 9999} \end{matrix}$$

### Step 3: Receive output

$$L \text{ shape} = (T \times D)$$

↑      ↙  
No. of words      Dim of embedding matrix

### Step 4: Training Model

- Dense vector can be used to train model
- During training, words are positioned with similar meanings together.  
Like "king" & "queen" vector values will be very close.

## Pretrained Word vectors

- precomputed embedded matrix trained on large text corpus
- Embedding layer weights fixed, only other layers are trained
- reduces training time & enhances accuracy (esp. if data is insufficient)

\* index to map words to vectors should not be '0'

'0' is used for padding

## Tensor Flow handling of text

Document  $\rightarrow$  token

1) Break down array of strings into single word strings

$\hookrightarrow$  known as Tokenization

$\hookrightarrow$  Each word in array is known as token.

$\hookrightarrow$  TensorFlow Tokenizer

$\hookrightarrow$  1) Break down string into tokens

2) convert list of tokens to list of Integers.

\* Tokenizer (Num\_words = K)

$\uparrow$   
only keep K most imp't words.

the rest assigned to <unknown> token.

2) Padding

$\hookrightarrow$  sequence of integers have different lengths (due to different sentence lengths)

$\hookrightarrow$  RNN do not accept different sequence lengths.

pad\_sequences() method.

padding = "pre" | "post"

$\hookrightarrow$  depending on task, padding location matters.

spam detection - want padding at start - can forget info.

language translation - want padding at end - need info at start to translate.

Data = N x T

$\uparrow$   $\uparrow$   
No. sentence (padded) sentence length

3) Embedding Layer

$\bullet$  N x T data passed through embedding layer to get corresponding dense vector

$\bullet$  Result is N x T x D

$\uparrow$   $\uparrow$   $\uparrow$  vector shape  
No. sentence length  
(samples) Each index have D sized vector  
Each sample have T indices