

LSTM - before

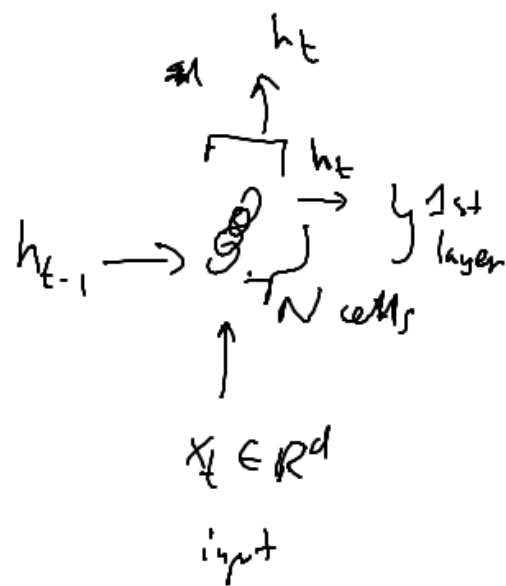
GRU - Gated Recurrent Unit (Cho et al., 2014)

$$z_t = \sigma_g \left(W_z \overset{N \times d \in \mathbb{R}^d}{x_t} + U_z \overset{N \times N \in \mathbb{R}^N}{h_{t-1}} + b_z \right) \quad \text{update gate}$$

$$r_t = \sigma_g \left(W_r x_t + U_r h_{t-1} + b_r \right) \quad \text{reset gate}$$

$$\hat{h}_t = \phi_h \left(W_h x_t + U_h (r_t \odot h_{t-1}) + b_h \right) \quad \text{candidate activation}$$

$$h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \hat{h}_t \quad \text{- output \& new hidden state}$$



$$\left\{ \begin{array}{l} \text{LSTM:} \\ c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t \end{array} \right.$$

$$\frac{\text{LSTM}}{4(Nd + N^2 + N)} \quad \text{weights}$$

$$\frac{\text{GRU}}{3(Nd + N^2 + N)} \quad \text{weights}$$

$$\frac{\text{simple RNN}}{Nd + N^2 + N}$$

Example

1 LSTM cell

input $(x_1, x_2, \dots, x_t) \in \mathbb{R}^t$

aim: predict $\hat{x}_{t+1} = x_1$

$(\underbrace{(x_1, x_2, \dots, x_t)}_{\text{randomly sampled}}, x_1)$



Natural Language Processing

Example: sentiment analysis

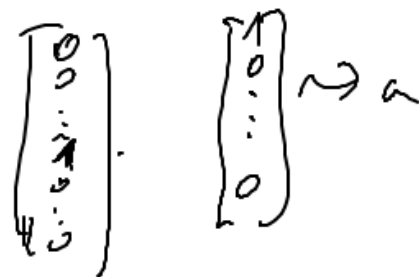
movie reviews $\xrightarrow{\text{aim}}$ predict how good was the movie according to the reviewer

'This movie is great!' \rightarrow 5/5
'This movie is horrible.' \rightarrow 1/5

one-hot encoding of words:

- 1) construct a dictionary of words, maybe with some special tokens e.g. [UNK] - unknown word token
[padding]
[eos]

- 2) each word is encoded as a 1-hot vector



Drawbacks:

- the dimension of the input is large
 - performance problems
 - training set needs to be very large

Optimal Solution: use word embeddings

i.e. a function

$$f: \{\text{words}\} \rightarrow \mathbb{R}^N$$

e.g. GloVe: # words = 400'000

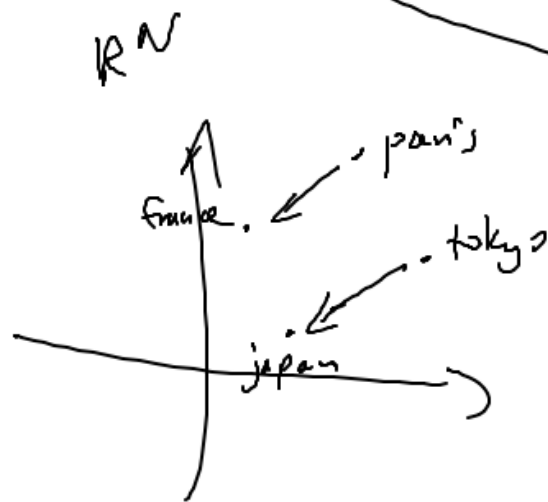
$N = 50, 100, 200, 300$

that can catch the meaning of words, for example
f should have similar values for words with similar meaning

In word2vec and GloVe the embedding function f has also this property; e.g.

$$f('paris') - f('france') \approx f('tokyo') - f('japan')$$

Words



similarity measure
usually are used:
(minus) Euclidean distance

$u \approx v$: ~~is~~
measure:

$$\frac{u \cdot v}{\|u\| \cdot \|v\|} \quad (= \cos \angle(u, v))$$

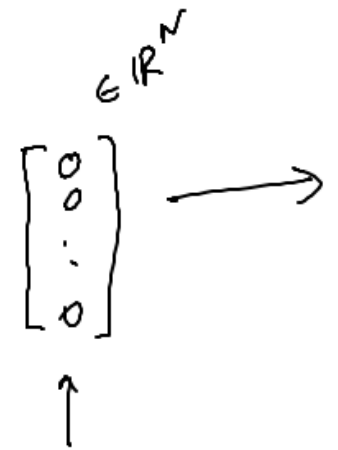
By using such embedding, the dim. of input is reduced and also we can have a smaller training set.

How to construct such an embedding?

Work 2vec



$W = \# \text{ words in the dict.}$



o_w one hot encoding of the word

$W \cdot o_w = e_w$

$\sigma(W \cdot o_w)$ \uparrow encoding of the word w

We feed a word c (context) and ~~also~~ check the prob. of word t (target) appearing close to c in the corpus.

$$\frac{\exp(V_j^T \cdot e_i)}{\sum_{j=1}^W \exp(V_j^T \cdot e_i)}$$

$$\mathcal{L}(y, \hat{y}) = - \sum_{j=1}^W y_j \cdot \log \hat{y}_j$$

E.g.

c - context word

t - target word chosen in a close proximity of c (say, ± 5 words)

one may want to sample the most common words (es) often,

the least - | - more often than we would

with a uniform distribution.

GloVe

construct a matrix of word-word occurrence probabilities (occurrence of cases where w_1 and w_2 are close): X_{ij} = probab. (frequency) of finding word i close to j

Prob. / ratio	k = solid	k = gas	k = water	k = fashion
$P(k ice)$	$1.9 \cdot 10^{-4}$	$6.6 \cdot 10^{-5}$	$3.0 \cdot 10^{-3}$	$1.7 \cdot 10^{-5}$
$P(k steam)$	$2.2 \cdot 10^{-5}$	$7.8 \cdot 10^{-4}$	$2.2 \cdot 10^{-3}$	$1.8 \cdot 10^{-5}$
$P(k ice)/P(k steam)$	8.9	$8.5 \cdot 10^{-2}$	1.36	0.96

$\theta_i \cdot e_j \approx \log X_{ij}$
 $\theta_i \cdot e_k \approx \log X_{ik}$
 $\theta_i \cdot (e_j - e_k) = \log \frac{X_{ij}}{X_{ik}}$

objective: find $\theta_i, e_i \in \mathbb{R}^N$, $i=1, \dots, W$ and b_i, b'_i , $i=1, \dots, W$

to minimize:

$$\sum_{i=1}^W \sum_{j=1}^W f(X_{ij}) \left(\theta_i \cdot e_j + b_i + b'_j - \log X_{ij} \right)^2$$

\uparrow weight, = 0 if $X_{ij} = 0$
 smaller for words which are very common, larger for less common

\rightarrow by product = embedding

smaller for words which are very common, larger for less common