

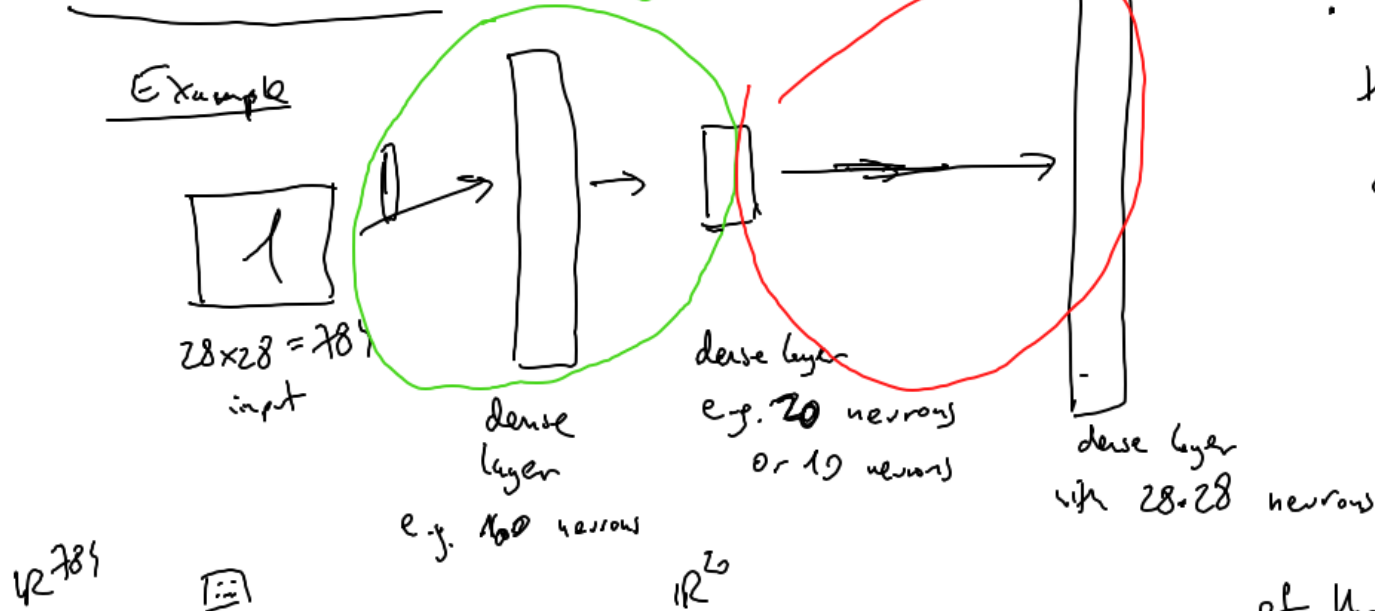
• LSTM (long-short term memory) layers 190

• GRU (gated recurrent unit) layers (Cho et al. 2014)
a simplified version of LSTM (smaller number of weights), but still offer the same performance

self-supervised learning

Auto encoders

Example



• We train the network so that the output is the same as the input

In fact, we want to obtain 2 networks which are parts of the above one:
encoder & decoder

Encoder reduces the dimension of the input without losing much information

The layers need not be dense, we could also have convolutional layers ~~and~~ etc.

Natural Language Processing (NLP)

Example: . sentiment analysis

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

'This movie is good' \rightarrow 4/5

'This movie is not good' \rightarrow 1/5

Cousy
 \updownarrow
Bad

. translations

. suggestion for next words (texting)

Good ... Monday Afternoon —

one-hot encoding of words:

- construct a dictionary of words, possibly with some special tokens

e.g. [UNK] unknown word

[padding] (if one wants to have sequences of fixed length)

[eos] end of sentence

- each word is encoded as a 1-hot vector

$\begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \rightarrow \text{home}$

$\begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \\ \vdots \\ 0 \end{bmatrix} \rightarrow a$

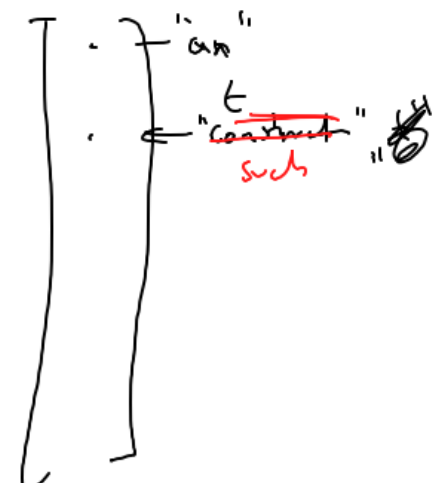
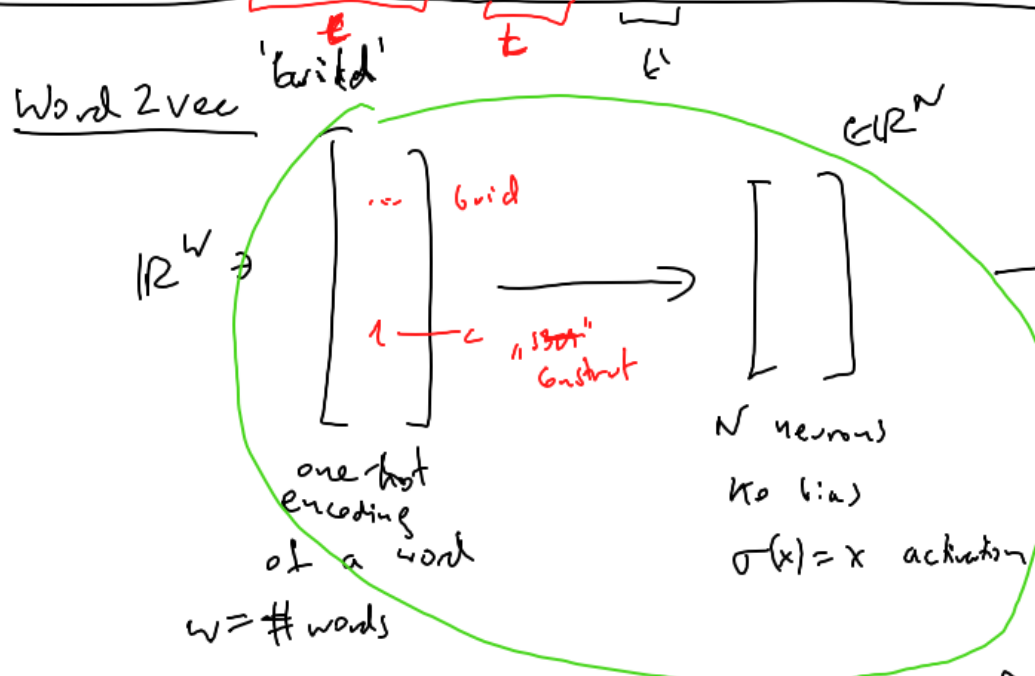
$\begin{bmatrix} 0 \\ 0 \\ 1 \\ \vdots \\ \vdots \end{bmatrix} \rightarrow \text{bad}$

$\begin{bmatrix} 0 \\ \vdots \\ 1 \end{bmatrix} \rightarrow \text{'busy'}$

Sequence of words \longrightarrow sequence of vectors (with 0 and 1)

The drawback: • the vectors have large dimension
• the training set has to be very large, so that the model 'sees' every word often enough

How to construct such an embedding?



c - word (context), input

t - target word: a word that is found close to the word c

This is not the real objective we care about this part

Glove

X_{ij} = probability (frequency) of finding word i close to j

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

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

j, j'

weight = 0 if $X_{ij} = 0$

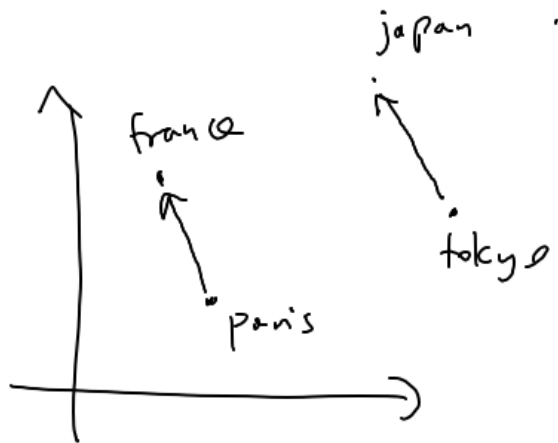
small for words that are very common (e.g. 'the')

	k=solid	k=gas	k=water	k=fashion
$P(k ice)$	$1.9 \cdot 10^{-4}$	$6.6 \cdot 10^{-5}$	$3.0 \cdot 10^3$	$1.7 \cdot 10^{-5}$
$P(k steam)$	$2.2 \cdot 10^{-5}$	$7.8 \cdot 10^{-4}$	$2.2 \cdot 10^3$	$1.8 \cdot 10^{-5}$
$\frac{P(k ice)}{P(k steam)}$	8.9	0.085	1.36	0.96

↓ embedding
as a by-product

These embeddings have some interesting properties:

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



it catches analogies

The dangers of using word embeddings:

Paper "Man is to Computer Programmer as
Woman is to Homemaker? Debiasing Word Embeddings"

~~Computer programmer~~ ... (A)
A man ... (B)



~~Man~~
Computer programmer ... (B)
A woman ... (A)

Jürgen ...
Janice ...
Mohammed ...

