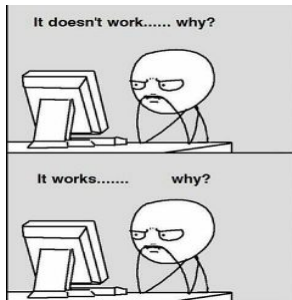


Lab class

Computational Semantics 2021

Adam Ek



The vocabulary

- A torchtext Field is
 - 1 a procedure of preprocessing steps
 - 2 a mapping of words to indices (numbers)
- Thus, we don't want a context Field and a center Field, but a shared Field.
- Since in fact, we only have one vocabulary

Projection function

```
TOKENS = bla bla bla

# read the csv files
train, test = TabularDataset.splits(path = ddir,
                                   train = 'train1.csv',
                                   test  = 'test1.csv',
                                   format = 'csv',
                                   fields = [('column1', TOKENS),
                                           ('column2', TOKENS)],
                                   skip_header = True,
                                   csv_reader_params = {'delimiter': '\t',
                                                       'quotechar': '½'})
```

What is an embedding?

- An embedding is:

$$e_{cat} = [x_1 x_2 \dots x_m] \quad (1)$$

- the individual numbers have *no meaning*
- embeddings obtain meaning through linear and non-linear transformations
- embeddings obtain meaning from what we can *predict* using them
- embeddings obtain meaning from their position in *space*

Model parameters

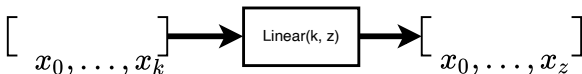
■ An embedding layer

$$\begin{bmatrix} x_0^0, \dots, x_k^0 \\ \vdots \\ x_0^V, \dots, x_k^V \end{bmatrix} \times \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix} = [x_0^2, \dots, x_k^2]$$

embedding layer one-hot word embedding

■ A prediction layer

Where z is the number of words in our vocabulary, and k the size of our vector



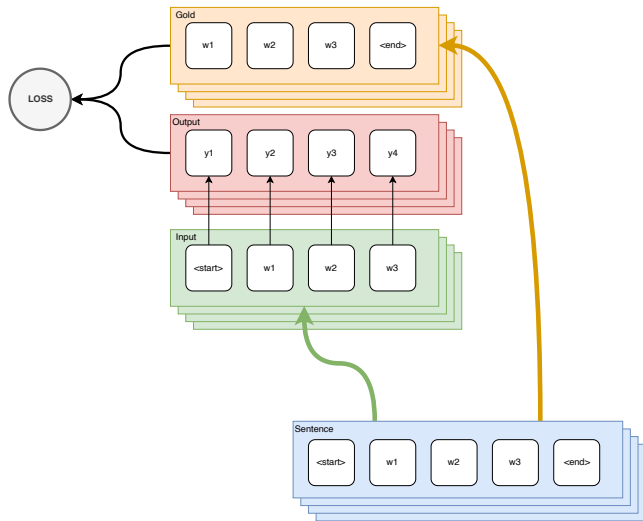
Projection function

Projection =	sum	sum	sum	sum	sum	sum
context word 1 =	6	1	4	3	1	1
context word 2 =	2	1	8	3	0	5
context word 3 =	4	1	2	3	9	5
New context =	12	3	14	9	10	11

Cross Entropy Loss

► Cross Entropy Loss

Language modeling



Open questions

