

—Neural modeling of language—

- What instruments for a connectionist modeling of language?
- What relationship with the semantic paradigm?
- distributional/vectoral for natural languages?

8a.1

1

▯ Semantic representation

q Principle of compositionality (recap)

- the meaning of an expression is function
 - the meaning of its subexpressions and • the way in which they are combined

q Cognitive motivation

- finite means (mind-brain) can deal with object potentially infinite (language)

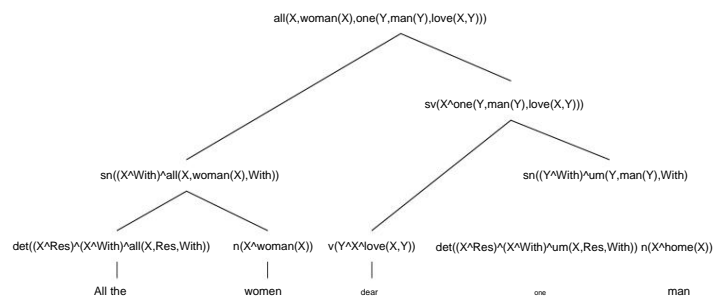
q Technological orientation

- assigning semantic representations to words, through rules appropriate for their combination, the representation is obtained semantics for any sentence...
- ...and from this representation we can build practical applications

8a.2

2

Inferential semantics: logical repr. (recap)



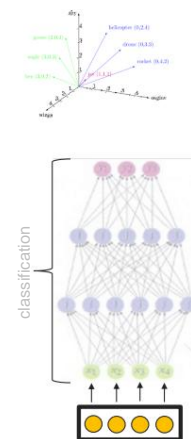
8a.3

3

Distributional semantics: vector repr. (recap)

q Lexical semantics

- meaning of the word represented by a vector
- vector (fixed size) of **high dimensionality** (eg 512)
- vector **condenses information** about the respective distribution, ie co-occurrence with other words



8a.4

4

Sentence and textual semantics 1/2

q What vector representation of the sentence/text?

• meaning of the text represented by vector: **what information** will be condensed by the respective vector?

• how to deal with arbitrary and very varied text lengths: **the same** vector size for any text length?

• compositionality principle: how to **combine** lexical vectors?

8a.5

5

Sentence and textual semantics 2/2

q What representation of the sentence/text?

neural architecture?

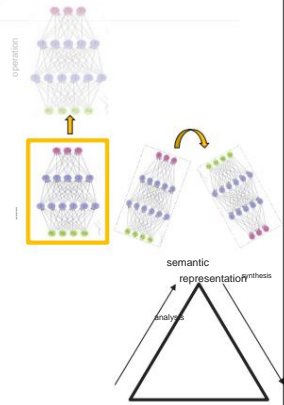
• how to obtain suitable representation as input to neural networks?

• A: with **other neural (sub-)networks!**

• how to get desired output?

• A: analogy with triangulation, but the training is joint ("**end-to-end**")

• with **what network architecture?** question motivating for this classroom



8a.6

6

Bag of words (BOW)

q Concatenation

$$\vec{v} = v_1 v_2 \dots = [v_{11}, \dots, v_{1m}, v_{21}, \dots, v_{2m}, \dots]$$

Just concatenate each vector i think

q Soma

$$\vec{v} = v_1 + v_2 + \dots =$$

$$[v_{11} + v_{21} + \dots, \dots, v_{1m} + v_{2m} + \dots]$$

q Average

You sum the word vectors component-wise, resulting in a single vector of the same size as a word vector.

$$\vec{v} = \text{aver}(v_1, v_2, \dots, v_n) =$$

$$[(v_{11} + v_{21} + \dots) / n, \dots, (v_{1n} + v_{2n} + \dots) / n]$$

q Maximum

Like the sum, but normalized by the number of words.

$$\vec{v} = \text{max}(v_1, v_2, \dots) =$$

$$[\max(v_{11}, v_{21}, \dots), \dots, \max(v_{1m}, v_{2m}, \dots)]$$

Take the maximum value in each dimension across all word vectors. This can emphasize salient features.

q "Zero degree" of

learning

• but in practice with interesting results, eg for similarity semantics, given the (little) sophistication required

8a.7

7

Convolution ("convolution")

q Rede Neuronal Convolutional ("CNN")

• the same network receives successive segments of *k* words

• another perspective: sliding window over *k*-grams

• for each *k*-gram it generates a vector

A CNN will look at small groups of words at a time, for example, 2 or 3 words at once.

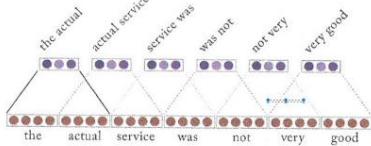
These small groups are called *k*-grams (*k* = number of words in the group).

q Aggregation ("pooling")

• vector representations of *k*-grams are aggregated into a vector

• max or average of each dimension of the vectors

We combine them into a single vector to represent the whole sentence.



8a.8

8

words are lexically ambiguous (recap) (recap)

WordNet Search - 3.1
- WordNet home page - Glossary - Help

Word to search for: kite Search WordNet

Display Options: (Select option to change) Change

Key: "S" = Show Synset (semantic) relations, "W" = Show Word (lexical) relations
Display options for sense: (gloss) "an example sentence"

Noun

- S₁ (n) kite (a bank check that has been fraudulently altered to increase its face value)
- S₂ (n) kite (a bank check drawn on insufficient funds at another bank in order to take advantage of the float)
- S₃ (n) kite (plaything consisting of a light frame covered with tissue paper; flown in wind at end of a string)
- S₄ (n) kite (any of several small graceful hawks of the family Accipitridae having long pointed wings and feeding on insects and small animals)

Verb

- S₁ (v) kite (increase the amount (of a check) fraudulently) "He kited many checks"
- S₂ (v) kite (get credit or money by using a bad check) "The businessman kited millions of dollars"
- S₃ (v) kite (soar or fly like a kite) "The pilot kited for a long time over the mountains"
- S₄ (v) kite (fly a kite) "Kids were kiting in the park"; "They kited the Red Dragon model"

8a.13

13

Recurrence

q Recurrent Neural Networks ("RNN")

- \tilde{y} internal state progressively compounded at each time step
- \tilde{y} at each time step, (the vector of) an input word and the state previous internal contribute to the current internal state and to the output
- \tilde{y} in the final step, the internal state is the representation (vector) of the sentence

8a.14

14

Compositionality

q Maintain memory of the past (during training)

- $\tilde{y} h_{t-1.V}$ contribution of the previous state
- $\tilde{y} x_{t.U}$ contribution of input at t to internal state
- $\tilde{y} h_t = g(h_{t-1.V} + x_{t.U})$ joint contribution
- $\tilde{y} h_{t.W}$ exit

q Model the language

- \tilde{y} V, U and W parameters learned / trained networks
- \tilde{y} training task: continuation task: *Peter said <cont?>*
- \tilde{y} learned model: language model (probabilistic)

8a.15

15

seq-to-seq

q Transduction (during inference)

- \tilde{y} eg translation, summarization, dialogue,...
- \tilde{y} final state of the encoder used as the basis for generating the output

q Encoder

- \tilde{y} RNN trained as a language model used to represent the meaning of the input

q Decoder ("decoder")

- \tilde{y} second RNN trained to generate output text
- \tilde{y} lexicon-sized output layer: next word predicted by the most active unit

8a.16

16

Resilience

q Transduction OK

• eg continuation, translation, summarization, dialogue

q Long-distance dependencies OK

• e.g, The *cats_i* sleep on the *pillows_k*. Mary told us that the Pedro *-as_k* washed it in the washing machine.

• same parameters/networks learned over each step
temporal/input word allows weighing and accommodating non-local relations

8a.17

17

Shortcomings 1/3

q Gradients

• learning/training involves many computations of the gradients and their factorization: with approximation values, errors accumulate

• creates bias to capture short-distance relationships

q Exploding gradients

• when many very large values > 1 occur • solution:
gradient clipping to lower the value

8a.18

18

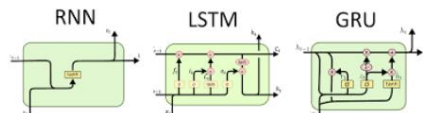
Shortcomings 2/3

q Exploding gradients q Vanishing gradients

• when many very small values < 1 occur

• solution: control the values that are brought from the past with gated cells:
most popular proposals (advanced topic):

• LSTM - Long Short Term Memory • GRU – Gated Recurrent Unit



8a.19

19

I think this is word orderer disambiguation, or something with the same word not being able to mean different things in different instances

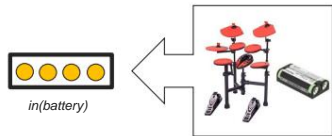
Shortcomings 3/3

q There is still no WSD

• undesirably, each word continues to have a single representation (same vector) for all its meanings/meanings: "blur" undefined semantic

• there is no word sense disambiguation (WSD: "word sense" disambiguation") depending on the context of occurrence

• e.g, Pedro bought an electric *drum* . To listen to it in the wireless headphones, without disturbing your neighbors, you just need to charge their *battery*



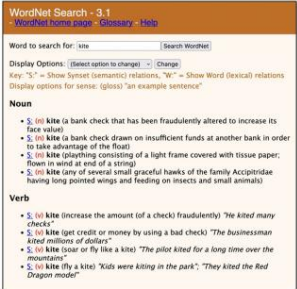
8a.20

20

Attention

Contextualized lexical disambiguation

- a word can have several possible meanings/meanings (in the lexicon)
- typically, one meaning is conveyed (in one occurrence)
- one of the meanings is favored, the others are ignored (by context)



8a.21

21

Attention

Contextualized lexical disambiguation

- a word can have several possible meanings/meanings (in the lexicon)
- typically, one meaning is conveyed (in one occurrence)
- one of the meanings is favored, the others are ignored (by context)

Alternative: Contextualized lexical representation

- semantic “blur” of a word is made more definite
- components of your vector are weighted according to the vectors of the other words in the context (and the task to be learned) • How? a solution in the following slides...

8a.22

22

dot product (recap)

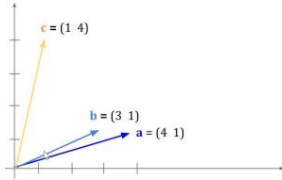
Scalar product of vectors (“dot product”)

$$z = v_1 \cdot v_2 = v_{11} \cdot v_{21} + \dots + v_{1m} \cdot v_{2m}$$

exs.:

$$a \cdot b = 13 = 3 \cdot 4 + 1 \cdot 1$$

$$a \cdot c = 8 = 1 \cdot 4 + 4 \cdot 1$$



highest result, closest vectors

recap: word embeddings: closer vectors, higher similarity semantics

8a.23

23

Attention to context 1/3

1. Obtain semantic representation

- vectorize (“embedding”): vectors (dimension m) for each of the n words $w_1 \dots w_n$ of a given sentence
- $v_1 = [v_{11}, \dots, v_{1m}]$
- ...
- $v_n = [v_{n1}, \dots, v_{nm}]$

2. Score semantic proximity

- scalar product of each word for each word (between vectors v_i) • $s_1 = [s_{11}, \dots, s_{1n}] = [v_1 \cdot v_1, \dots, v_1 \cdot v_n]$
- ...
- $s_n = [s_{n1}, \dots, s_{nn}]$

8a.24

24

Attention to context 2/3

3. Obtain weights with softmax normalization

• **normalize**: each component lies within [0,1] and sums to 1

• $\tilde{y} \text{ wg1} = \text{softmax}(s1) = [\text{wg11}, \dots, \text{wg1n}]$

• $\tilde{y} \dots$

• $\tilde{y} \text{ wgn} = \text{softmax}(sn) = [\text{wgn1}, \dots, \text{wgnn}]$

4. Add heavy starting vectors

• \tilde{y} contextualized representation of each wdi: sum of the vectors of departure of the wdi words in the context after being weighed by the weight vector of each wdi

• $\tilde{y} \text{ y1} = \text{wg11.v1} + \dots + \text{wg1n.vn}$

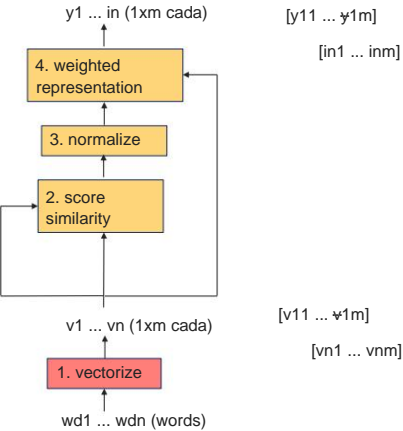
• $\tilde{y} \dots$

• $\tilde{y} \text{ in} = \text{wgn1.v1} + \dots + \text{wgnn.vn}$

8a.25

25

Attention to context 3/3



8a.26

26

Attention to task 1/3

q OK: Context-sensitive weighting $w_1 \dots w_n$

• based on the dot product

• lexical vectors

• terminology: self-attention: only input information taken into consideration

q BUT ALSO: Task-based weighting

• eg summarization, translation,...

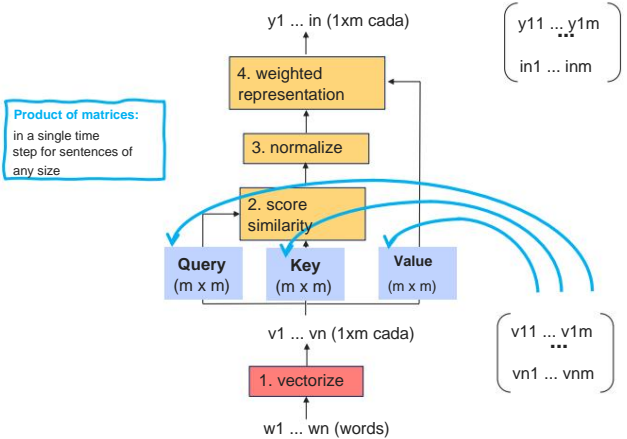
• based on the matrix product

• matrices/network weights learned during training for the task

8a.27

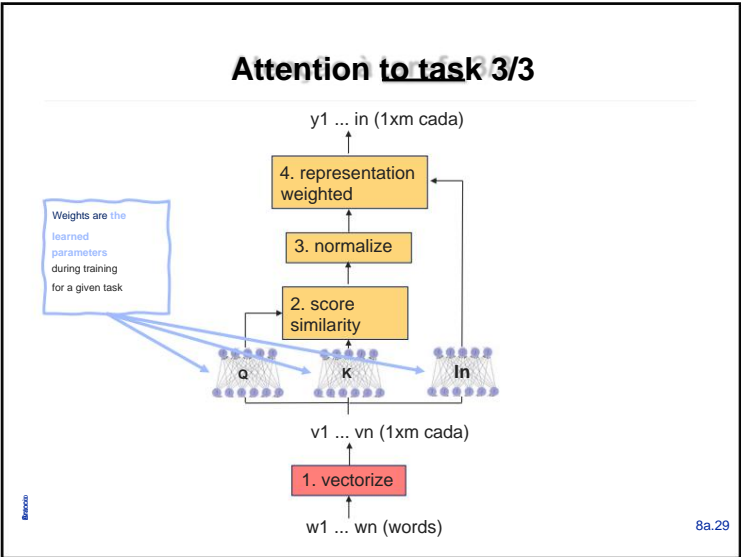
27

Attention to task 2/3

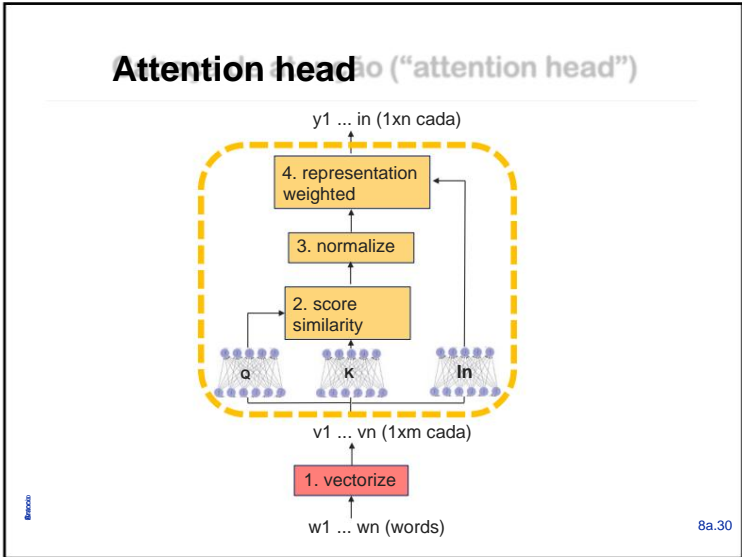


8a.28

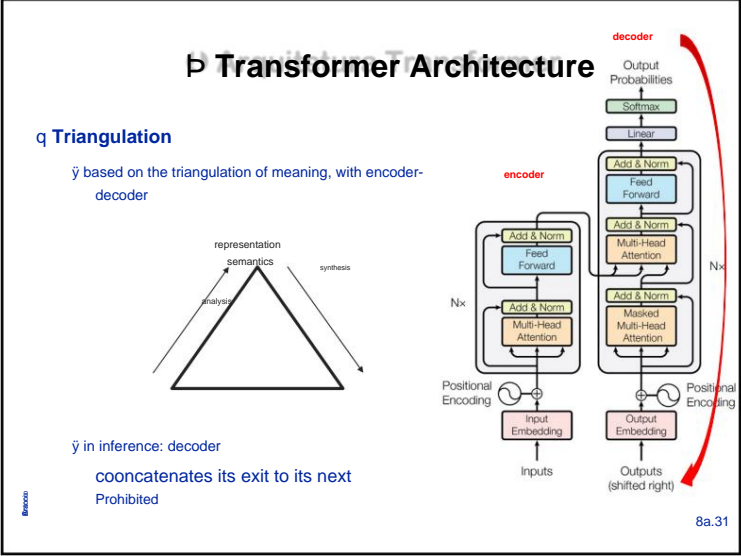
28



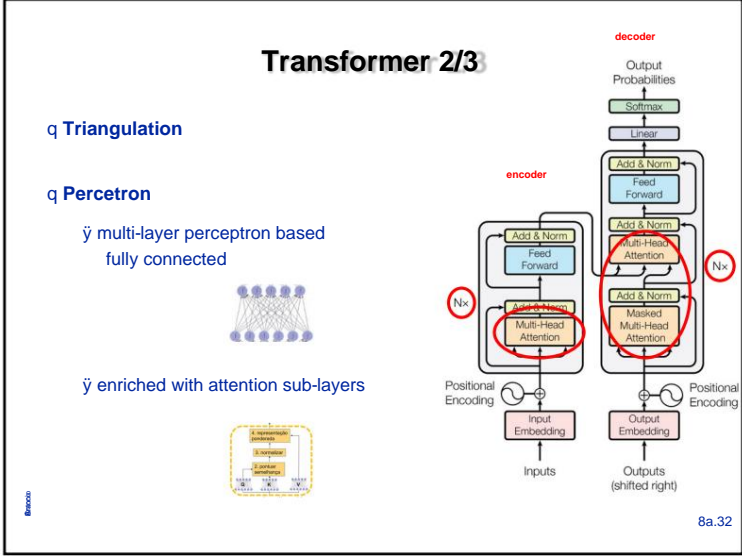
29



30



31



32

Transformer 3/3

q Triangulation

q Percetron

q Disambiguation

ȳ based on "defined" semantic representations, i.e. contextualized, i.e. reciprocally disambiguated

8a.33

33

Key innovations

q parallelism I

ȳ input words processed in parallel, in a single self-attention time step (GPU)

ȳ in training: codif/desc; in test: codif

ȳ Overcoming explosive gradients of RNNs

q parallelism II

ȳ contextualized word representations do not need to be collapsed and diluted into a semantic representation of the sentence (before the last output layer)

ȳ overcoming the problem of RNNs in long-range forgetting

8a.34

34

Key adaptations

q relative word order

ȳ reassured with vectors that represent the positions in the sentence ("positional embeddings"), to be added to the word vectors

q transduction

ȳ parallel decoder trained with masked: for each word, hide the words after that

what attention

ȳ no decoder: V receives representation of previously output words, Q and K receive representation of all input words by the encoder

8a.35

35

Technicalities 1/2 /2

q Multi-head attention

ȳ At each layer, the input (set of words) is passed through several heads in parallel, outputs are concatenated and linearized

ȳ each head learns differently regularities and linguistic relations of the context and the task

8a.36

36

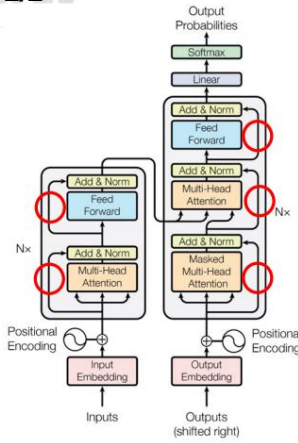
Technicalities 2/2

q Residual connections

- bypass each sub-layer to add your input to your output
- facilitate backpropagation during training

q Some hyperparameters

- (no original paper from 2017)
- lexical vectors size 512
- 6 layers
- 8 heads per layer (size 64)



8a.37

37

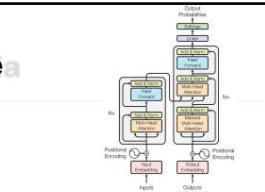
Reference

q Reference

- Vaswani et al., 2017, *Attention Is All You Need*, Google Research

q Mainstream

- Transformer is the mainstream architecture, with best performance for a wide range of NLP tasks
- Transformer variants: next lessons



8a.38

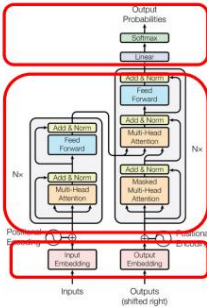
38

- Conclusion -

q Index

- Vector representation of texts and semantic triangulation
- Neural network architectures: CNN, RNN
- Attention
- Transformer Architecture

How to get the input word vectors?
How to get the output from the last layer?
How to train a model?



8a.39

39