Compositionality in Semantic Spaces

Martha Lewis

ILLC University of Amsterdam

2nd Symposium on Compositional Structures

Glasgow, UK

1 Categorical Compositional Distributional Semantics

2 Shifting Categories

3 Recursive Neural Networks



1 Categorical Compositional Distributional Semantics

2 Shifting Categories

- 3 Recursive Neural Networks
- 4 Summary and Outlook

When a male octopus spots a female, his normally grayish body suddenly becomes striped. He swims above the female and begins caressing her with seven of his arms.

Cherries jubilee on a white suit? Wine on an altar cloth? Apply club soda immediately. It works beautifully to remove the stains from fabrics.

Steven Pinker. The Language Instinct: How the Mind Creates Language (Penguin Science) (pp. 1-2).

When a male octopus spots a female, his normally grayish body suddenly becomes striped. He swims above the female and begins caressing her with seven of his arms.

Cherries jubilee on a white suit? Wine on an altar cloth? Apply club soda immediately. It works beautifully to remove the stains from fabrics.

Steven Pinker. The Language Instinct: How the Mind Creates Language (Penguin Science) (pp. 1-2).

... And how can we get computers to do the same?

Compositional Distributional Semantics



Compositional Distributional Semantics



Distributional hypothesis

Words that occur in similar contexts have similar meanings [Harris, 1958].

?

?

?

?

?

?

?

?

?

?

?

?

?

?

?

- U.S. Senate, because they are It made him
 - It made him
- sympathy for the problems of
- peace and the sanctity of
- without the accompaniment of
- a monstrous crime against the
- this mystic bond between the
- suggests a current nostalgia for
 - Harbor" in 1915), the
 - an earthy and very
 - To be
 - Ordinarily, the
 - nothing in the whole range of
 - It is said that fear in
 - megatons: the damage to

- , like to eat as high on the
- beings caught up in the life are not only religious sacrifice.
- race.
- and natural world that the
- values in art.
- element was the compelling modern dance work,
- , he believes, is to seek one's liver synthesizes only enough
- liver synthesizes only enough
- experience more widely
- beings produces an odor that
- germ plasm would be such

Distributional hypothesis

Words that occur in similar contexts have similar meanings [Harris, 1958].

U.S. Senate, because they are human , like to eat as high on the It made him human sympathy for the problems of human beings caught up in the peace and the sanctity of human life are not only religious without the accompaniment of human sacrifice. a monstrous crime against the human race. this mystic bond between the human and natural world that the suggests a current nostalgia for human values in art. Harbor" in 1915), the **human** element was the compelling an earthy and very human modern dance work, To be human , he believes, is to seek one's Ordinarily, the **human** liver synthesizes only enough nothing in the whole range of **human** experience more widely It is said that fear in human beings produces an odor that megatons: the damage to human germ plasm would be such

Distributional Semantics

- Words are represented as vectors
- Entries of the vector are derived from how often the target word co-occurs with the context word



Similarity is given by cosine distance:

$$sim(v,w) = \cos(heta_{v,w}) = rac{\langle v,w
angle}{||v||||w||}$$

The role of compositionality

Compositional distributional models

We can produce a sentence vector by composing the vectors of the words in that sentence.

$$\overrightarrow{s} = f(\overrightarrow{w_1}, \overrightarrow{w_2}, \dots, \overrightarrow{w_n})$$

Three generic classes of CDMs:

- Vector mixture models [Mitchell and Lapata (2010)]
- *Tensor-based* models [Coecke, Sadrzadeh, Clark (2010); Baroni and Zamparelli (2010)]
- Neural models [Socher et al. (2012); Kalchbrenner et al. (2014)]

Applications (1/2)

Why are CDMs important?

The problem of producing robust representations for the meaning of phrases and sentences is at the heart of every task related to natural language.

• Paraphrase detection

Problem: Given two sentences, decide if they say the same thing in different words

Solution: Measure the cosine similarity between the sentence vectors

• Sentiment analysis

Problem: Extract the general sentiment from a sentence or a document

Solution: Train a classifier using sentence vectors as input

Textual entailment

Problem: Decide if one sentence logically infers a different one Solution: Examine the feature inclusion properties of the sentence vectors

Machine translation

Problem: Automatically translate one sentence into a different language

Solution: Encode the source sentence into a vector, then use this vector to decode a surface form into the target language

• And so on. Many other potential applications exist...

A general programme

- 1. a Choose a compositional structure, such as a pregroup or combinatory categorial grammar.
 - b Interpret this structure as a category, the grammar category.
- 2. a Choose or craft appropriate meaning or concept spaces, such as vector spaces, density matrices, or conceptual spaces.
 - b Organize these spaces into a category, the **semantics category**, with the same abstract structure as the grammar category.
- 3. Interpret the compositional structure of the grammar category in the semantics category via a functor preserving the necessary structure.
- 4. Bingo! This functor maps type reductions in the grammar category onto algorithms for composing meanings in the semantics category.



Quantizing the grammar

Coecke, Sadrzadeh and Clark (2010):

Pregroup grammars are structurally homomorphic with the category of finite-dimensional vector spaces and linear maps (both share compact closure)

• In abstract terms, there exists a structure-preserving passage from grammar to meaning:

 $\mathcal{F}:\mathsf{Grammar}\to\mathsf{Meaning}$

The meaning of a sentence w₁w₂...w_n with grammatical derivation α is defined as:

$$\overrightarrow{w_1w_2\ldots w_n} := \mathcal{F}(\alpha)(\overrightarrow{w_1}\otimes \overrightarrow{w_2}\otimes \ldots \otimes \overrightarrow{w_n})$$

A pregroup grammar $P(\Sigma, B)$ is a relation that assigns grammatical types from a Compact CC freely generated over a set of atomic types B to words of a vocabulary Σ .

• Atomic types $x \in \mathcal{B}$ have morphisms

$$\begin{split} \epsilon_x^r &: x \cdot x^r \to 1, \qquad \epsilon_x^l : x^l \cdot x \to 1 \\ \eta_x^r &: 1 \to x^r \cdot x, \qquad \eta_x^l : 1 \to x \cdot x^l \end{split}$$

- Elements of the pregroup are basic (atomic) grammatical types, e.g. B = {n, s}.
- Atomic grammatical types can be combined to form types of higher order (e.g. n · n^l or n^r · s · n^l)
- A sentence $w_1 w_2 \dots w_n$ (with word w_i to be of type t_i) is grammatical whenever:

$$t_1 \cdot t_2 \cdot \ldots \cdot t_n \rightarrow s$$

Pregroup derivation: example

$$p \cdot p^{r} \rightarrow 1 \rightarrow p^{r} \cdot p \qquad p^{l} \cdot p \rightarrow 1 \rightarrow p \cdot p^{l}$$

$$Sad clowns tell jokes$$

$$n \cdot n^{l} \cdot n \cdot n^{r} \cdot s \cdot n^{l} \cdot n \rightarrow n \cdot 1 \cdot n^{r} \cdot s \cdot 1$$

$$= n \cdot n^{r} \cdot s$$

$$\rightarrow 1 \cdot s$$

$$= s$$

We define a strongly monoidal functor $\ensuremath{\mathcal{F}}$ such that:

 $\mathcal{F}: P(\Sigma, \mathcal{B}) \to \textbf{FVect}$

$$\begin{array}{rcl} \mathcal{F}(p) &=& P \quad \forall p \in \mathcal{B} \\ \mathcal{F}(1) &=& \mathbb{R} \\ \mathcal{F}(p \cdot q) &=& \mathcal{F}(p) \otimes \mathcal{F}(q) \\ \mathcal{F}(p^{r}) = \mathcal{F}(p^{l}) &=& \mathcal{F}(p) \\ \mathcal{F}(p \leq q) &=& \mathcal{F}(p) \rightarrow \mathcal{F}(q) \\ \mathcal{F}(\epsilon^{r}) = \mathcal{F}(\epsilon^{l}) &=& \text{inner product in FVect} \\ \mathcal{F}(\eta^{r}) = \mathcal{F}(\eta^{l}) &=& \text{identity maps in FVect} \end{array}$$

[Kartsaklis, Sadrzadeh, Pulman and Coecke, 2016]

The grammatical type of a word defines the vector space in which the word lives:

- Nouns are vectors in N;
- adjectives are linear maps $N \rightarrow N$, i.e elements in $N \otimes N$;
- intransitive verbs are linear maps N → S, i.e. elements in N ⊗ S;

 transitive verbs are bi-linear maps N ⊗ N → S, i.e. elements of N ⊗ S ⊗ N;

• The composition operation is tensor contraction, i.e. elimination of matching dimensions by application of inner product.

Graphical calculus: example



 $\mathcal{F}(\alpha)(\overline{\text{trembling}} \otimes \overline{\text{shadows}} \otimes \overline{\text{play}} \otimes \overline{\text{hide-and-seek}})$



- Formal semantic approaches are good at composition...
- ... but the things they compose are featureless atoms
- Distributional semantics give a much richer meaning to their atoms...
- ... but have no obvious compositional mechanisms.
- The categorical compositional model marries the two by interpreting both the grammar and the sematics category as being compact closed.
- The structure and interactions of the grammar category are mapped over to the semantic category.

Categorical Compositional Distributional Semantics

2 Shifting Categories

3 Recursive Neural Networks



A general programme

- a Choose a compositional structure, such as a pregroup or combinatory categorial grammar.
 - b Interpret this structure as a category, the grammar category.
- 2. a Choose or craft appropriate meaning or concept spaces, such as vector spaces, density matrices, or conceptual spaces.
 - b Organize these spaces into a category, the **semantics category**, with the same abstract structure as the grammar category.
- 3. Interpret the compositional structure of the grammar category in the semantics category via a functor preserving the necessary structure.
- 4. Bingo! This functor maps type reductions in the grammar category onto algorithms for composing meanings in the semantics category.

Conceptual spaces [Gärdenfors, 2014] can provide a more cognitively realistic semantics.



$\textit{noun} \in \textit{COLOUR} \otimes \textit{SHAPE} \otimes \cdots$

Convex algebras

- Notation. For a set X we write $\sum_{i} p_i |x_i\rangle$ for a finite formal convex sum of elements of X, where $p_i \in \mathbb{R}^{\geq 0}$ and $\sum_{i} p_i = 1$.
- A convex algebra is a set A with a mixing operation α satisfying:

$$\alpha(|a\rangle) = a$$

$$\alpha(\sum_{i,j} p_i q_{i,j} |a_{i,j}\rangle) = \alpha(\sum_i p_i |\alpha(\sum_j q_{i,j} |a_{i,j}\rangle)))$$

- Examples: Real or complex vector spaces, simplices, semilattices, trees
- A convex relation between convex algebras (A, α) and (B, β) is a relation that commutes with forming convex combinations, i.e.

$$(\forall i.R(a_i, b_i)) \Rightarrow R(\sum_i p_i a_i, \sum_i p_i b_i)$$

- We define the category **ConvexRel** as having convex algebras as objects and convex relations as morphisms
- **ConvexRel** is compact closed with \times as monoidal product.
- We build a functor from pregroup grammar in the same way: choose a space *N* for nouns, a space *S* for sentences.
- Drawback: how do we start to build word representations?

- We can use density matrices, and more generally, positive operators rather than vectors to represent words.
- Positive operators A, B have the Löwner ordering $A \sqsubseteq B \iff B A$ is positive.
- We use this ordering to represent entailment, and introduce a graded version useful for linguistic phenomena.
- We will show that graded entailment lifts compositionally to sentence level.
- Similar approaches in [Sadrzadeh et al., 2018]

Words as positive operators

- A positive operator P is self-adjoint and $orall v \in V. \left< v \right| P \left| v \right> \geq 0$
- Density matrices are convex combinations of projectors: $\rho = \sum_{i} p_{i} |v_{i}\rangle \langle v_{i}|, \text{ s.t. } \sum_{i} p_{i} = 1$
- We can view concepts as collections of items, with *p_i* indicating relative frequency.
- For example:

$$egin{aligned} &\left[\!\left[pet
ight]\!\right] =& p_d \left| dog
ight
angle \left\langle dog
ight| + p_c \left| cat
ight
angle \left\langle cat
ight| + \ & p_t \left| tarantula
ight
angle \left\langle tarantula
ight| + ... \end{aligned}$$
 where $orall i.p_i \geq 0$ and $\sum_i p_i = 1$

 There are various choices for normalisation of the density matrices...

Sentence meaning in **CPM(FVect**)

- We assign semantics via a strong monoidal functor S : $\textbf{Preg} \rightarrow \textbf{CPM}(\textbf{FVect})$
- Let $w_1 w_2 ... w_n$ be a string of words with corresponding grammatical types t_i in $\mathbf{Preg}_{\{n,s\}}$ such that $t_1, ... t_n \xrightarrow{r} s$
- Let $[\![w_i]\!]$ be the meaning of word w_i in **CPM(FVect)**. Then the meaning of $w_1 w_2 \dots w_n$ is given by:

$$\llbracket w_1 w_2 \dots w_n \rrbracket = \mathsf{S}(r)(\llbracket w_1 \rrbracket \otimes \dots \otimes \llbracket w_n \rrbracket)$$



So how do we do graded hyponymy?

- Recall that positive operators A, B have the Löwner ordering $A \sqsubseteq B \iff B A$ is positive.
- We say that A is a hyponym of B if $A \sqsubseteq B$
- We say that A is a k-hyponym of B for a given value of k in the range (0,1] and write A ≼_k B if:

B - kA is positive

• We are interested in the maximum such k.

Theorem

For positive self-adjoint matrices A, B such that $supp(A) \subseteq supp(B)$, the maximum k such that $B - kA \ge 0$ is given by $1/\lambda$ where λ is the maximum eigenvalue of B^+A .

k-hyponymy interacts well with compositionality

- We would like our notion of entailment to work at the sentence level.
- Since sentences are represented as positive operators, we can compare them directly.
- If sentences have similar structure, we can also give a lower bound on the entailment strength between sentences based on the entailment strengths between the words in the sentences.

Example

Suppose $\llbracket dog \rrbracket \preccurlyeq_k \llbracket pet \rrbracket$ and $\llbracket park \rrbracket \preccurlyeq_l \llbracket field \rrbracket$. Then

 $\llbracket My \text{ dog runs in the park} \rrbracket \preccurlyeq_{???} \llbracket My \text{ pet runs in the field} \rrbracket$

How should we build density matrices for words?





WordNet. Princeton University. 2010

M. Lewis	Semantic Spaces	30/53	
----------	-----------------	-------	--

Hyperlex [Vulić et al., 2017] gold-standard dataset. 2,163 noun pairs human annotated.

Model	Dev	Test
Poincaré embeddings	-	0.512
SDSN (Random)	0.757	0.692
SDSN (Lexical)	0.577	0.544
Density matrices	-	0.551
Density matrices (non)	-	0.631

Dataset from [Sadrzadeh et al., 2018] consisting of 23 sentence pairs. Example pairs are:

 $recommend \ development \models suggest \ improvement \\ progress \ reduce \models \ development \ replace$

With normalization:

Model	ρ	р
Verb-only	0.268	> 0.25
Frobenius mult.	0.508	> 0.05
Frobenius n.c.	0.436	> 0.05
Additive	0.643	> 0.001
Inter-annotator	0.66	-

[Bankova et al., 2019]

Dataset from [Sadrzadeh et al., 2018]. Example pairs are: recommend development \models suggest improvement progress reduce \models development replace

Without normalization:

Model	ρ	р
Verb-only	0.370	> 0.1
Frobenius mult.	0.696	$> 5\cdot 10^{-4}$
Frobenius n.c.	0.306	0.15
Additive	0.737	$>5\cdot10^{-5}$
Inter-annotator	0.66	-



Figure 1: Learned diagonal variances, as used in evaluation (Section 6), for each word, with the first letter of each word indicating the position of its mean. We project onto generalized eigenvectors between the mixture means and variance of query word *Bach*. Nearby words to *Bach* are other composers e.g. *Mozart*, which lead to similar pictures.

34/53

[Vilnis and McCallum, 2014], and similar approaches seen in [Jameel and Schockaert, 2017], [Bražinskas et al., 2017]

- Vector spaces provide a good means of talking about similarity between words.
- But both conceptual spaces and positive operators give us a richer word representation.
- We don't have a good way of building word representations in conceptual spaces yet, but using information from WordNet to build density matrices seems to give good results with very little effort.

Categorical Compositional Distributional Semantics

2 Shifting Categories

- 3 Recursive Neural Networks
- 4 Summary and Outlook

Recursive Neural Networks



A (compact closed) category for neural networks

???

Alternatively... linear recursive neural networks



[Lewis, 2019]

Alternatively... linear recursive neural networks



Alternatively... linear recursive neural networks



Adjectives and intransitive verbs





Content/structure split





himself : ns^r n^{rr} n^r s



- Understanding neural networks as a semantics category for Lambek categorial grammar opens up more possibilities to use tools from formal semantics in computational linguistics.
- We can immediately see possibilities for building alternative networks.
- Decomposing the tensors for functional words into repeated applications of a compositionality function gives options for learning representations.
- Brittleness of word types in DisCo is alleviated

- Incorporate non-linearity
- Extend to other types of network that are currently state-of-the-art
- Testing on data comparison with standard DisCo representations, examine ability to switch word types, look at specialised tensors

Categorical Compositional Distributional Semantics

2 Shifting Categories

3 Recursive Neural Networks



- The categorical compositional model of [Coecke et al., 2010] can be instantiated with a choice of grammar category (which we saw yesterday in Sadrzadeh's and Wijnhold's talks).
- We also have a choice of meaning category, allowing us to move towards a richer semantics. Initial results are positive.
- We have started to link to neural network methods for building word vectors.

- Develop ways of building word regions from corpora
- Links between positive matrices and multivariate Gaussian word embeddings?
- Building on more cognitively plausible concept representations.
- Starting to model meaning change and game-theoretic models of language [Bradley et al., 2018, Hedges and Lewis, 2018].
- Modelling non-literal uses of language.

NWO Veni grant 'Metaphorical Meanings for Artificial Agents'

References I



Bankova, D., Coecke, B., Lewis, M., and Marsden, D. (2019).

Graded entailment for compositional distributional semantics. arXiv preprint arXiv:1601.04908. to appear in Journal of Language Modelling.



Bowman, S. R., Potts, C., and Manning, C. D. (2014).

Recursive neural networks can learn logical semantics. arXiv preprint arXiv:1406.1827.



Bradley, T.-D., Lewis, M., Master, J., and Theilman, B. (2018).

Translating and evolving: Towards a model of language change in discocat. arXiv preprint arXiv:1811.11041.



Bražinskas, A., Havrylov, S., and Titov, I. (2017).

Embedding words as distributions with a bayesian skip-gram model. *arXiv preprint arXiv:1711.11027*.



Coecke, B., Sadrzadeh, M., and Clark, S. (2010).

Mathematical Foundations for a Compositional Distributional Model of Meaning. Lambek Festschrift. *Linguistic Analysis*, 36:345–384.



Gärdenfors, P. (2014).

The geometry of meaning: Semantics based on conceptual spaces. MIT Press.



Hedges, J. and Lewis, M. (2018).

Towards functorial language-games. arXiv preprint arXiv:1807.07828.



References II



Jameel, S. and Schockaert, S. (2017).

Modeling context words as regions: An ordinal regression approach to word embedding. In Proceedings of the 21st Conference on Computational Natural Language Learning (CoNLL 2017), pages 123–133.



Lewis, M. (2019).

Compositionality for recursive neural networks. *IfCoLog Journal of Applied Logics*. to appear.



Moulton, D., Goriely, A., and Chirat, R. (2015). The morpho-mechanical basis of ammonite form. *Journal of Theoretical Biology*, 364:220–230.



Sadrzadeh, M., Kartsaklis, D., and Balkır, E. (2018). Sentence entailment in compositional distributional semantics.

Annals of Mathematics and Artificial Intelligence, 82(4):189–218.



Socher, R., Huval, B., Manning, C., and A., N. (2012). Semantic compositionality through recursive matrix-vector spaces. In *Conference on Empirical Methods in Natural Language Processing 2012*.



Vilnis, L. and McCallum, A. (2014).

Word representations via gaussian embedding. arXiv preprint arXiv:1412.6623.



Vulić, I., Gerz, D., Kiela, D., Hill, F., and Korhonen, A. (2017). Hyperlex: A large-scale evaluation of graded lexical entailment. *Computational Linguistics*, 43(4):781–835.

