Modeling Language and Translations, Categorically

Tai-Danae Bradley CUNY Graduate Center

Martha Lewis University of Amsterdam Jade Master UC Riverside Brad Theilman UC San Diego

syntax

semantics

syntax

semantics

The Yoneda lemma for linguistics:

"You shall know a word by the company it keeps."

– John Firth, 1957

(a.k.a. the distributional hypothesis)



syntax semantics

students
$$\mapsto \begin{bmatrix} 8 \\ 5 \\ 0 \end{bmatrix} = \text{students} \in \mathbb{R}^3$$

where \mathbb{R}^3 is generated by

$$\mathbf{school} = \begin{bmatrix} 1\\0\\0 \end{bmatrix} \quad \mathbf{books} = \begin{bmatrix} 0\\1\\0 \end{bmatrix} \quad \mathbf{butterfly} = \begin{bmatrix} 0\\0\\1 \end{bmatrix}$$

syntax semantics $\|$ (FVect, \otimes , \mathbb{R})

We choose our semantics category to be the category of **finite-dimensional real vector spaces**. It's a symmetric monoidal category. It's also compact closed (every object has a dual).

- monoidal product = tensor product
- monoidal unit = ground field

semantics syntax cordially adverb ۶. students learn verb noun in preposition purple adjective sweater noun

syntax semantics $(\mathsf{G},\cdot,1)$

We choose our syntax category to be the free compact closed category $(G, \cdot, 1)$ on a finite set of grammar types. Example: if the set is $\{n, s\}$:



Definition: Let G be a free compact closed category on a finite set of grammar types. A distributional categorical language model—or **language model**, for short—is a strong monoidal functor

$$(\mathsf{G},\cdot,1) \xrightarrow{F} (\mathsf{FVect},\otimes,\mathbb{R})$$

Every grammar type g in G corresponds to a vector space Fgand every grammar reduction $r: g \to h$ gives rise to a linear transformation $Fr: Fg \to Fh$. (pause)

A little background:

- In 2010, Coecke et. al. model language via the Cartesian product G × FVect.
 - Coecke, B., Sadrzadeh M., and Clark, S. "Mathematical foundations for a compositional distributional model of meaning." *arXiv:1003.4394*
- In 2014, Kartsaklis et. al. model language via a strong monoidal functor $G \rightarrow FVect_V$.
 - Kartsaklis, D., Sadrzadeh, M., Pulman, S. and Coecke, B. "Reasoning about meaning in natural language with compact closed categories and Frobenius algebras." *Logic and Algebraic Structures in Quantum Computing and Information*, p. 199; *arXiv:1401.5980*

(unpause)

Definition: Let G be a category that is freely monoidal on a finite set of grammar types. A distributional categorical language model—or **language model**, for short—is a strong monoidal functor

$$(\mathsf{G},\cdot,1) \xrightarrow{F} (\mathsf{FVect},\otimes,\mathbb{R})$$

To pair grammar types with meaning vectors, e.g.

$$n \rightsquigarrow \begin{bmatrix} 8 & 5 & 0 \end{bmatrix}^{\top} = \mathbf{students}^{\top}$$

we wish to use the Grothendieck construction.

But... F is not Cat-valued!



where CV is the category whose objects are vectors $\mathbf{v} \in V$ and there is a unique morphism $\mathbf{v} \to \mathbf{v}'$ labeled by the Euclidean distance $d(\mathbf{v}, \mathbf{v}')$. For any linear map $f: V \to W$ define $Cf: CV \to CW$ to be the functor that agrees with fon objects and is given by $d(\mathbf{v}, \mathbf{v}') \mapsto d(f\mathbf{v}, f\mathbf{v}')$ on morphisms. **Definition**: Let F be a language model and let C be as before.

$$\mathsf{G} \xrightarrow{F} \mathsf{FVect} \xrightarrow{C} \mathsf{Cat}$$

The **product space representation** of F with respect to C, denoted $\int_C F := \int CF$, is the Grothendieck construction of CF. Explicitly, it is a category with

objects:morphisms: (g, \mathbf{w}) $(g, \mathbf{w}) \xrightarrow{(r,d)} (h, \mathbf{v})$ wherewhere $\mathbf{w} \in CFg$ $q \xrightarrow{r} h$ $d := d(CFr\mathbf{w}, \mathbf{v})$

Definition: Let F be a language model and let C be as before.

$$\mathsf{G} \xrightarrow{F} \mathsf{FVect} \xrightarrow{C} \mathsf{Cat}$$

The **product space representation** of F with respect to C, denoted $\int_C F := \int CF$, is the Grothendieck construction of CF. Explicitly, it is a category with

objects:

(n, students)

where

morphisms:

$$(g, \mathbf{w}) \xrightarrow{(r,d)} (h, \mathbf{v})$$

where

students $\in CFn$

$$g \xrightarrow{r} h \qquad d := d(CFr\mathbf{w}, \mathbf{v})$$

Proposition: Let F be a language model. Then $\int_C F$ is a monoidal category.

• On objects, the monoidal product is

$$(g, \mathbf{w}) \otimes (h, \mathbf{v}) = (gh, \Phi_{g,h}(\mathbf{w} \otimes \mathbf{v}))$$

• On morphisms, the monoidal product is

$$(r,d)\otimes(r',d')=(rr',\Phi_{g,h}(d\otimes d'))$$

where $\Phi_{g,h}: CFg \otimes CFg \rightarrow CF(gh)$ is the natural isomorphism included in the data of the strong monoidal functor CF.

Next, we want to make sense of the assignment

students $\mapsto (n, \text{students})$

Definition: Let F be a language model and let W be the free monoid on a finite set of words, viewed as a discrete category. A **lexicon** is a functor

$$\begin{array}{c} W \\ \downarrow \ell \\ \int_C F \end{array}$$

Next, we want to make sense of the assignment

students \mapsto (*n*, students)

Definition: Let F be a language model and let W be the free monoid on a finite set of words, viewed as a discrete category. A **lexicon** is a functor



Next, we want to make sense of the assignment

students $\mapsto (n, \text{students})$

Definition: Let F be a language model and let W be the free monoid on a finite set of words, viewed as a discrete category. A **lexicon** is a functor



We can model language by

- considering a monoidal functor grammar \rightarrow vector spaces
- use the Grothedieck construction to pair grammar types with meaning vectors
- view words in a text as a discrete category, and map them to the Grothendieck construction.



G' is free on $\{n_S, s_S\}$

G is free on $\{n_E, s_E\}$

Definition: A translation $T = (j, \alpha)$ from a language model F to a language model F' is a monoidal functor j and a monoidal natural transformation $\alpha \colon F \Rightarrow F'j$.



Definition: A translation $T = (j, \alpha)$ from a language model F to a language model F' is a monoidal functor j and a monoidal natural transformation $\alpha \colon F \Rightarrow F'j$.



$$\alpha_{n_E} \colon F(n_E) \to F'(n_s)$$

 $\alpha_{s_E} \colon F(s_E) \to F'(s_s)$

Example

$$j(n_E) = n_S$$

 $j(s_E) = s_S$
 $F(n_E) = N_E$
 $F'(n_S) = N_S$
 $F(s_E) = S_E$
 $F('s_S) = S_S$

DisCoCat

is the category with

objects:

morphisms:

language models

translations

$$(\mathsf{G},\cdot,1) \xrightarrow{F} (\mathsf{FVect},\otimes,\mathbb{R})$$

Distributional Compositional Categorical language models



Proposition. Let K: FVect \rightarrow Cat be a fully faithful functor and let MonCat denote the category of monoidal categories and strong monoidal functors. There is a functor



where $\int_K T$ is the strong monoidal functor given by...

(take K = C)





in Fh



Definition: Let $\ell: W \to \int_C F$ and $\ell': W' \to \int_C F'$ be lexicons, and let $T = (j, \alpha)$ be a translation from F to F'. The F-F' **dictionary** with respect to T is the comma category

$$\left(\int_C T \circ \ell\right) \downarrow \ell$$

Concretely, it is the set (discrete category) of triples (w, (r, d), w') where $w \in W$ and $w' \in W'$ and $(r, d) \colon \int_C T(\ell w) \to \ell' w'$ is a morphism in $\int_C F'$.

Example: Suppose W and W' are as below, and we have a translation $T = (j, \alpha)$ between language models F and F' (as before), where α can be determined appropriately.

$$W = \{ \text{students, learn, ...} \}$$

 $W' = \{ \text{estudiantes, aprenden, ...} \}$

then (students learn, (r, d), estudiantes aprenden) is an object in the F-F' dictionary, where $j(n_E n_E^r s_E) \xrightarrow{r} n_S n_S^r s_S$ is a reduction in G' and the distance d is determined by first translating $\mathbf{st} \otimes \mathbf{lrn}$ then applying the linear map corresponding to $F'r(\alpha_{nn^rs}(\mathbf{st} \otimes \mathbf{lrn}))$ the reduction.

in $F(n_s) \otimes F(n_s) \otimes F(s_S)$

Future Work

- Adapt the model to handle change in word order: e.g. "red car" vs. "coche rojo." (In this talk, we only considered the fragment of language consisting of nouns/intransitive verbs)
- 2. Take advantage of string diagram calculus.
- 3. Used an enriched version (Lawvere metric spaces).
- 4. Investigate meaning change and negotiated meaning between speakers (language evolution).

Thank you!