

A corpus-based toy model for DisCoCat

Stefano Gogioso

Quantum Group, University of Oxford

SLPCS 2016

The Plan

What we assume:

¹Free and finite-dimensional.

What we assume:

- (i) an abstract corpus, as a set/sequence of sentences

¹Free and finite-dimensional.

The Plan

What we assume:

- (i) an abstract corpus, as a set/sequence of sentences
- (ii) each sentence annotated with a constituent structure tree
 - we consider context-free grammars à la Chomsky.

¹Free and finite-dimensional.

The Plan

What we assume:

- (i) an abstract corpus, as a set/sequence of sentences
- (ii) each sentence annotated with a constituent structure tree
 - we consider context-free grammars à la Chomsky.

What we do with it:

¹Free and finite-dimensional.

The Plan

What we assume:

- (i) an abstract corpus, as a set/sequence of sentences
- (ii) each sentence annotated with a constituent structure tree
 - we consider context-free grammars à la Chomsky.

What we do with it:

- (i) obtain a toy pregroup grammar from the annotated corpus
 - entirely object-oriented, no sentence type

¹Free and finite-dimensional.

What we assume:

- (i) an abstract corpus, as a set/sequence of sentences
- (ii) each sentence annotated with a constituent structure tree
 - we consider context-free grammars à la Chomsky.

What we do with it:

- (i) obtain a toy pregroup grammar from the annotated corpus
 - entirely object-oriented, no sentence type
- (ii) obtain semantics in a category of R -semimodules¹
 - any involutive commutative semiring R , but here we focus on \mathbb{N}

¹Free and finite-dimensional.

The Plan

What we assume:

- (i) an abstract corpus, as a set/sequence of sentences
- (ii) each sentence annotated with a constituent structure tree
 - we consider context-free grammars à la Chomsky.

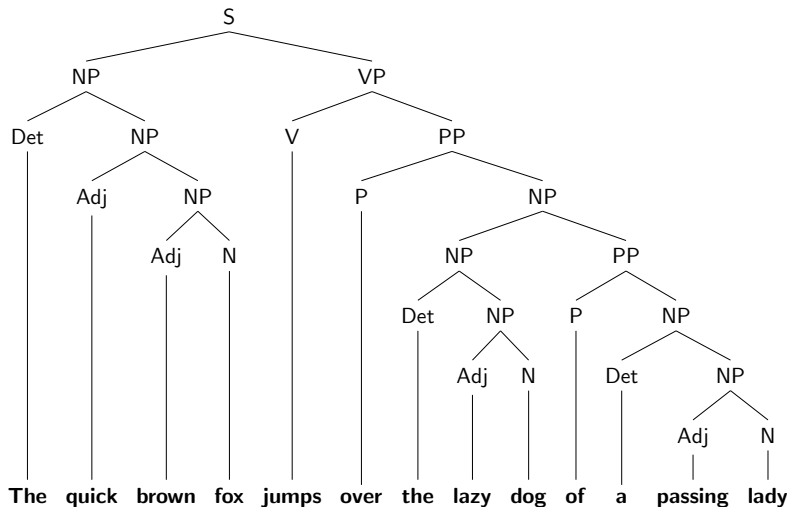
What we do with it:

- (i) obtain a toy pregroup grammar from the annotated corpus
 - entirely object-oriented, no sentence type
- (ii) obtain semantics in a category of R -semimodules¹
 - any involutive commutative semiring R , but here we focus on \mathbb{N}

Our semantics are free/minimal, in a certain sense explained later.

¹Free and finite-dimensional.

Constituent Structure Trees



A single atomic type n for *objects*, together with:

A single atomic type n for *objects*, together with:

- (i) object words, of type n
 - the nouns in the corpus sentences

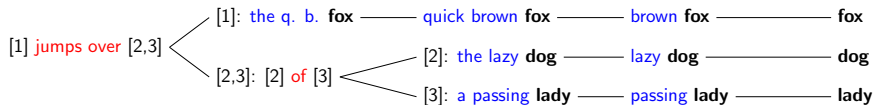
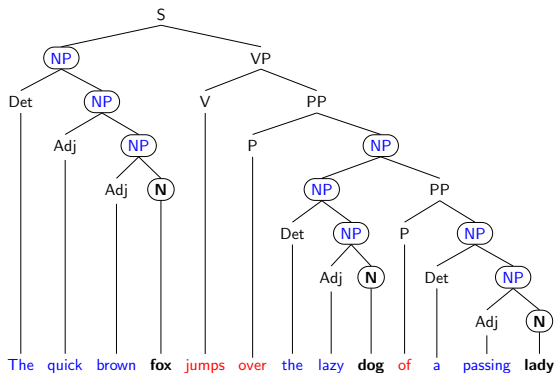
A single atomic type n for *objects*, together with:

- (i) object words, of type n
 - the nouns in the corpus sentences
- (ii) modifying words, of type $n^r \cdot n$ or $n \cdot n^l$
 - the words modifying nouns/NPs into other NPs

A single atomic type n for *objects*, together with:

- (i) object words, of type n
 - the nouns in the corpus sentences
- (ii) modifying words, of type $n^r \cdot n$ or $n \cdot n^l$
 - the words modifying nouns/NPs into other NPs
- (iii) interaction fragments, of type $n^r \cdot n \cdot n^l$
 - sentence fragments connecting noun phrases

Objects and their Interactions



Interaction words

Modifying words

Interaction fragments

The Category of R -Semimodules

To obtain our semantics, we consider:

²Including probability distributions and stochastic maps.

The Category of R -Semimodules

To obtain our semantics, we consider:

- (i) an involutive, commutative semiring R

²Including probability distributions and stochastic maps.

The Category of R -Semimodules

To obtain our semantics, we consider:

- (i) an involutive, commutative semiring R
- (ii) the category $R\text{-Mod}$ of free finite-dim R -semimodules
 - objects in the form R^X , for finite sets X

²Including probability distributions and stochastic maps.

The Category of R -Semimodules

To obtain our semantics, we consider:

- (i) an involutive, commutative semiring R
- (ii) the category $R\text{-Mod}$ of free finite-dim R -semimodules
 - objects in the form R^X , for finite sets X
 - morphisms $R^X \rightarrow R^Y$ are $Y \times X$ R -valued matrices

²Including probability distributions and stochastic maps.

The Category of R -Semimodules

To obtain our semantics, we consider:

- (i) an involutive, commutative semiring R
- (ii) the category $R\text{-Mod}$ of free finite-dim R -semimodules
 - objects in the form R^X , for finite sets X
 - morphisms $R^X \rightarrow R^Y$ are $Y \times X$ R -valued matrices

Many examples of interest are in this form:

²Including probability distributions and stochastic maps.

The Category of R -Semimodules

To obtain our semantics, we consider:

- (i) an involutive, commutative semiring R
- (ii) the category $R\text{-Mod}$ of free finite-dim R -semimodules
 - objects in the form R^X , for finite sets X
 - morphisms $R^X \rightarrow R^Y$ are $Y \times X$ R -valued matrices

Many examples of interest are in this form:

- finite sets and relations, for $R = \mathbb{Bool}$

²Including probability distributions and stochastic maps.

The Category of R -Semimodules

To obtain our semantics, we consider:

- (i) an involutive, commutative semiring R
- (ii) the category $R\text{-Mod}$ of free finite-dim R -semimodules
 - objects in the form R^X , for finite sets X
 - morphisms $R^X \rightarrow R^Y$ are $Y \times X$ R -valued matrices

Many examples of interest are in this form:

- finite sets and relations, for $R = \mathbb{Bool}$
- finite-dim real/complex vector spaces, for $R = \mathbb{R}, \mathbb{C}$

²Including probability distributions and stochastic maps.

The Category of R -Semimodules

To obtain our semantics, we consider:

- (i) an involutive, commutative semiring R
- (ii) the category $R\text{-Mod}$ of free finite-dim R -semimodules
 - objects in the form R^X , for finite sets X
 - morphisms $R^X \rightarrow R^Y$ are $Y \times X$ R -valued matrices

Many examples of interest are in this form:

- finite sets and relations, for $R = \mathbb{Bool}$
- finite-dim real/complex vector spaces, for $R = \mathbb{R}, \mathbb{C}$
- finite-dim convex cones², for $R = \mathbb{R}^+$

²Including probability distributions and stochastic maps.

The Category of R -Semimodules

To obtain our semantics, we consider:

- (i) an involutive, commutative semiring R
- (ii) the category $R\text{-Mod}$ of free finite-dim R -semimodules
 - objects in the form R^X , for finite sets X
 - morphisms $R^X \rightarrow R^Y$ are $Y \times X$ R -valued matrices

Many examples of interest are in this form:

- finite sets and relations, for $R = \mathbb{Bool}$
- finite-dim real/complex vector spaces, for $R = \mathbb{R}, \mathbb{C}$
- finite-dim convex cones², for $R = \mathbb{R}^+$
- finite multi-sets and “multi-relations”, for $R = \mathbb{N}$

²Including probability distributions and stochastic maps.

The Category of R -Semimodules

Some desirable features of the category of R -semimodules:

The Category of R -Semimodules

Some desirable features of the category of R -semimodules:

- (i) $R\text{-Mod}$ is a \dagger -symmetric monoidal category

The Category of R -Semimodules

Some desirable features of the category of R -semimodules:

- (i) $R\text{-Mod}$ is a \dagger -symmetric monoidal category
- (ii) $R\text{-Mod}$ is compact closed, with self-dual objects

The Category of R -Semimodules

Some desirable features of the category of R -semimodules:

- (i) $R\text{-Mod}$ is a \dagger -symmetric monoidal category
- (ii) $R\text{-Mod}$ is compact closed, with self-dual objects
- (iii) $R\text{-Mod}$ has classical structures associated to canonical bases

$$\begin{array}{c} \text{---} \\ | \\ \circ \\ / \quad \backslash \\ \text{---} \quad \text{---} \end{array} := \sum_{x \in X} \begin{array}{c} \triangleup_x \quad \triangleleft_x \quad \triangleup_x \\ | \quad | \quad | \\ \text{---} \quad \text{---} \quad \text{---} \end{array} \qquad \begin{array}{c} \text{---} \\ | \\ \circ \\ | \\ \text{---} \end{array} := \sum_{x \in X} \begin{array}{c} \triangleleft_x \\ | \\ \text{---} \end{array}$$

The Choice of $R = \mathbb{N}$

Why would we want to choose $R = \mathbb{N}$?

The Choice of $R = \mathbb{N}$

Why would we want to choose $R = \mathbb{N}$?

- (i) vectors are covariant over morphisms $f : R \rightarrow S$ of semirings
- e.g. vector $(v_x)_{x \in X} \in R^X$ is mapped to $(f(v_x))_{x \in X} \in S^X$

The Choice of $R = \mathbb{N}$

Why would we want to choose $R = \mathbb{N}$?

- (i) vectors are covariant over morphisms $f : R \rightarrow S$ of semirings
 - e.g. vector $(v_x)_{x \in X} \in R^X$ is mapped to $(f(v_x))_{x \in X} \in S^X$
- (ii) \mathbb{N} is initial in the category of commutative semirings

The Choice of $R = \mathbb{N}$

Why would we want to choose $R = \mathbb{N}$?

- (i) vectors are covariant over morphisms $f : R \rightarrow S$ of semirings
 - e.g. vector $(v_x)_{x \in X} \in R^X$ is mapped to $(f(v_x))_{x \in X} \in S^X$
- (ii) \mathbb{N} is initial in the category of commutative semirings
 - \Rightarrow the choice of \mathbb{N} is the least restrictive possible

The Choice of $R = \mathbb{N}$

Why would we want to choose $R = \mathbb{N}$?

- (i) vectors are covariant over morphisms $f : R \rightarrow S$ of semirings
 - e.g. vector $(v_x)_{x \in X} \in R^X$ is mapped to $(f(v_x))_{x \in X} \in S^X$
- (ii) \mathbb{N} is initial in the category of commutative semirings
 - \Rightarrow the choice of \mathbb{N} is the least restrictive possible

Semantics in any other semiring R of interest can be recovered by covariance over the unique morphism $\mathbb{N} \rightarrow R$.

The Distributional Part

We construct our semantic space from the corpus:

The Distributional Part

We construct our semantic space from the corpus:

- (i) consider the set X of all word instances in all sentences

$$X = \{(\underline{s}, j) \mid \underline{s} \text{ sentence, } j \text{ index of word instance in } \underline{s}\}$$

The Distributional Part

We construct our semantic space from the corpus:

- (i) consider the set X of all word instances in all sentences

$$X = \{(\underline{s}, j) \mid \underline{s} \text{ sentence, } j \text{ index of word instance in } \underline{s}\}$$

- (ii) take R^X as the semantic space

The Distributional Part

We construct our semantic space from the corpus:

- (i) consider the set X of all word instances in all sentences

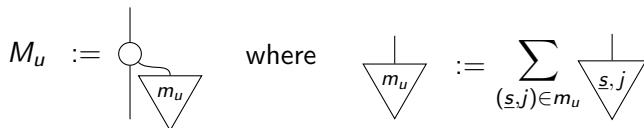
$$X = \{(\underline{s}, j) \mid \underline{s} \text{ sentence, } j \text{ index of word instance in } \underline{s}\}$$

- (ii) take R^X as the semantic space
- (iii) embed a word w as the indicator function of all its instances

$$\begin{array}{c} \downarrow \\ \triangle \\ w \end{array} := \sum_{s_j=w} \begin{array}{c} \downarrow \\ \triangle \\ \underline{s}, j \end{array}$$

The Compositional Part

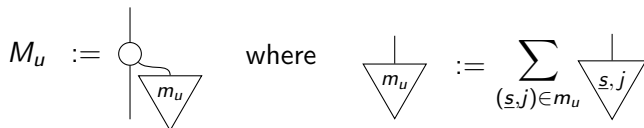
Modifier words are mapped to projectors:

$$M_u := \text{circle} \text{---} \text{triangle}_{m_u} \quad \text{where} \quad \text{triangle}_{m_u} := \sum_{(\underline{s}, j) \in m_u} \text{triangle}_{\underline{s}, j}$$


³As long as R satisfies the additive cancellation law.

The Compositional Part

Modifier words are mapped to projectors:

$$M_u := \text{circle} \text{---} \text{triangle}(m_u) \quad \text{where} \quad \text{triangle}(m_u) := \sum_{(\underline{s}, j) \in m_u} \text{triangle}(\underline{s}, j)$$


We define m_u to be the set of instances of object words which appear in objects modified by u . For example, from our sentence we'd have:

$$\begin{aligned} \{\text{fox}\} &\subseteq m_{\text{quick}} \\ \{\text{lady}\} &\subseteq m_{\text{passing}} \\ \{\text{fox}, \text{dog}\} &\subseteq m_{\text{the}} \end{aligned}$$

³As long as R satisfies the additive cancellation law.

The Compositional Part

Modifier words are mapped to projectors:

$$M_u := \begin{array}{c} | \\ \circ \\ | \end{array} \quad \text{where} \quad \begin{array}{c} | \\ \nabla \\ m_u \end{array} := \sum_{(\underline{s}, j) \in m_u} \begin{array}{c} | \\ \nabla \\ \underline{s}, j \end{array}$$

We define m_u to be the set of instances of object words which appear in objects modified by u . For example, from our sentence we'd have:

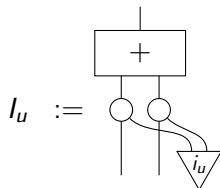
$$\begin{aligned} \{\text{fox}\} &\subseteq m_{\text{quick}} \\ \{\text{lady}\} &\subseteq m_{\text{passing}} \\ \{\text{fox}, \text{dog}\} &\subseteq m_{\text{the}} \end{aligned}$$

We automatically get some logic out of operator algebra³.

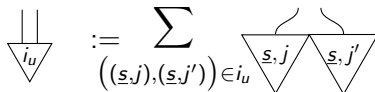
³As long as R satisfies the additive cancellation law.

The Compositional Part

Interaction fragments are mapped to binary operations. First we construct the operation for single-word fragments:



where



The Compositional Part

Interaction fragments are mapped to binary operations. First we construct the operation for single-word fragments:

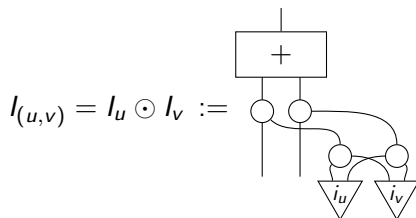


We define i_u to be the set of pairs of instances which appear in objects put into relation by u . For example, from our sentence we'd have:

$$\{(fox,dog), (fox,lady)\} \subseteq i_{jumps} \quad (0.1)$$

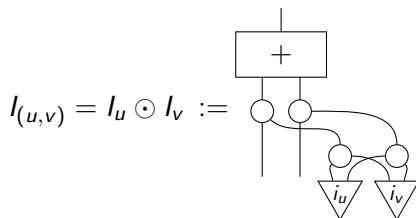
The Compositional Part

We use Frobenius algebras to treat multi-word fragments:



The Compositional Part

We use Frobenius algebras to treat multi-word fragments:

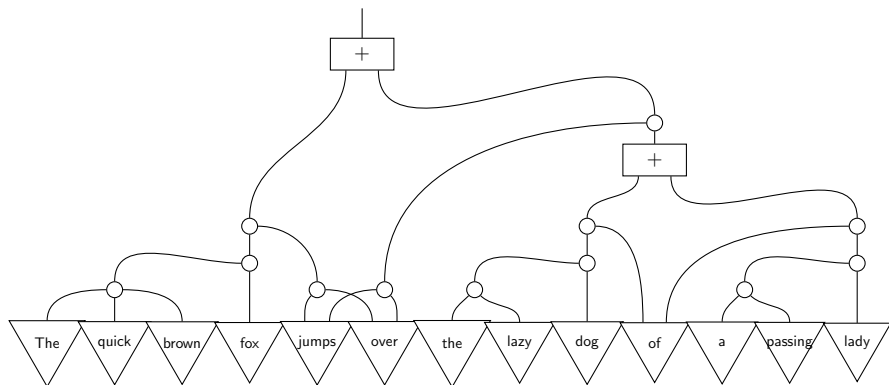


The linear operator $\boxplus : R^X \otimes R^X \rightarrow R^X$ is defined as follows:

$$\boxplus := \left(id_{R^X} \otimes \eta + \eta \otimes id_{R^X} \right) = |\underline{s}, j\rangle \otimes |\underline{s}', j'\rangle \mapsto |\underline{s}, j\rangle + |\underline{s}', j'\rangle$$

The End Result

Here is the resulting⁴ semantics for our sentence:



⁴After some applications of the spider theorem, to group modifier words together.

Future work

A lot of things to do!

A lot of things to do!

- (i) Toy model needs a number of improvements
 - treatment of personal/possessive pronouns
 - treatment of conjunctions

A lot of things to do!

- (i) Toy model needs a number of improvements
 - treatment of personal/possessive pronouns
 - treatment of conjunctions
- (ii) More sophisticated choice of semiring
 - encoding of polarity, modality and inflection

A lot of things to do!

- (i) Toy model needs a number of improvements
 - treatment of personal/possessive pronouns
 - treatment of conjunctions
- (ii) More sophisticated choice of semiring
 - encoding of polarity, modality and inflection
- (iii) Compressing the free model to obtain concrete models
 - change of semiring + linear compression \Rightarrow more semantics?

A lot of things to do!

- (i) Toy model needs a number of improvements
 - treatment of personal/possessive pronouns
 - treatment of conjunctions
- (ii) More sophisticated choice of semiring
 - encoding of polarity, modality and inflection
- (iii) Compressing the free model to obtain concrete models
 - change of semiring + linear compression \Rightarrow more semantics?
- (iv) CPM construction (possibly iterated)
 - treatment ambiguity and entailment

A lot of things to do!

- (i) Toy model needs a number of improvements
 - treatment of personal/possessive pronouns
 - treatment of conjunctions
- (ii) More sophisticated choice of semiring
 - encoding of polarity, modality and inflection
- (iii) Compressing the free model to obtain concrete models
 - change of semiring + linear compression \Rightarrow more semantics?
- (iv) CPM construction (possibly iterated)
 - treatment ambiguity and entailment
- (v) Enriched/higher order categories
 - encode simplicial structure extracted from the corpus

Thanks for Your Attention!

Any Questions?