

Abstract Categorical Grammar Parsing

the general case

in Honor of Gérard Huet

Philippe de Groote
Inria-Lorraine

Content

- 1 Definition of ACG
- 2 Examples
- 3 Some Key Properties
- 4 Constructing a Parsing Algorithm

Definition

Motivations

- To provide a type-theoretic notion of grammar, taking advantages of ideas by Curry and Lambek.
- To provide a grammatical framework in which other existing grammatical models may be encoded.
- To see the parse-structures as first-class citizen.
- To allow the user to define grammatical composition combinators.
- To base the formalism on a small set of mathematical primitives that combine via simple composition rules.

Types, signatures and λ -terms

$\mathcal{T}(A)$ is the set of linear implicative types built on the set of atomic types A :

$$\mathcal{T}(A) ::= A \mid (\mathcal{T}(A) \multimap \mathcal{T}(A))$$

A higher-order linear signature is a triple $\Sigma = \langle A, C, \tau \rangle$, where:

A is a finite set of atomic types;

C is a finite set of constants;

$\tau : C \rightarrow \mathcal{T}(A)$ is a function that assigns each constant in C with a linear implicative type built on A .

$\Lambda(\Sigma)$ denotes the set of linear λ -terms built upon a higher-order linear signature Σ .

Vocabularies and Lexicons

A vocabulary is simply defined to be a higher-order linear signature.

Given two vocabularies $\Sigma_1 = \langle A_1, C_1, \tau_1 \rangle$ and $\Sigma_2 = \langle A_2, C_2, \tau_2 \rangle$, a lexicon $\mathcal{L} = \langle \eta, \theta \rangle$ from Σ_1 to Σ_2 is made of two functions:

$$\eta : A_1 \rightarrow \mathcal{T}(A_2),$$

$$\theta : C_1 \rightarrow \Lambda(\Sigma_2),$$

such that

$$\vdash_{\Sigma_2} \theta(c) : \eta(\tau_1(c)).$$

Definition

An abstract categorial grammar is a quadruple

$$\mathcal{G} = \langle \Sigma_1, \Sigma_2, \mathcal{L}, s \rangle$$

where :

$\Sigma_1 = \langle A_1, C_1, \tau_1 \rangle$ and $\Sigma_2 = \langle A_2, C_2, \tau_2 \rangle$ are two higher-order linear signatures; Σ_1 is called the abstract vocabulary and Σ_2 is called the object vocabulary;

$\mathcal{L} : \Sigma_1 \rightarrow \Sigma_2$ is a lexicon from the abstract vocabulary to the object vocabulary;

$s \in \mathcal{T}(A_1)$ is a type of the abstract vocabulary; it is called the distinguished type of the grammar.

Languages generated by an ACG

The abstract language generated by \mathcal{G} ($\mathcal{A}(\mathcal{G})$) is defined as follows:

$$\mathcal{A}(\mathcal{G}) = \{t \in \Lambda(\Sigma_1) \mid \vdash_{\Sigma_1} t: s \text{ is derivable}\}$$

The object language generated by \mathcal{G} ($\mathcal{O}(\mathcal{G})$) is defined to be the image of the abstract language by the term homomorphism induced by the lexicon \mathcal{L} :

$$\mathcal{O}(\mathcal{G}) = \{t \in \Lambda(\Sigma_2) \mid \exists u \in \mathcal{A}(\mathcal{G}). t = \mathcal{L}(u)\}$$

Some properties

- Membership is decidable if and only if Multiplicative Exponential Linear Logic is decidable.
- Membership for lexicalized ACGs is NP-complete.
- Membership for second-order ACGs is polynomial.

Examples

Strings as linear λ -terms

There is a canonical way of representing strings as linear λ -terms. It consists of representing strings as function composition:

$$/abbac/ = \lambda x. a (b (b (a (c x))))$$

In this setting:

$$\begin{aligned} \epsilon &\triangleq \lambda x. x \\ \alpha + \beta &\triangleq \lambda x. \alpha (\beta x) \end{aligned}$$

Signatures

Σ_0 : N, NP, S : type
 J : NP
 U : N
 A : $N \multimap ((NP \multimap S) \multimap S)$
 S : $((NP \multimap S) \multimap S) \multimap (NP \multimap S)$

Σ_1 : a, John, seeks, unicorn : *STRING*

Σ_2 :
 ι, o : type
 \wedge : $o \multimap (o \multimap o)$
 \exists : $(\iota \rightarrow o) \multimap o$
 j : ι
unicorn : $\iota \multimap o$
find : $\iota \multimap (\iota \multimap o)$
try : $\iota \multimap ((\iota \multimap o) \multimap o)$

Lexicons

$\mathcal{L}_1 : \Sigma_0 \rightarrow \Sigma_1$

$N, NP, S := STRING$
 $J := /John/$
 $U := /unicorn/$
 $A := \lambda x. \lambda p. p (/a/ + x)$
 $S := \lambda p. \lambda x. p (\lambda y. x + /seeks/ + y)$

$\mathcal{L}_2 : \Sigma_0 \rightarrow \Sigma_2$

$N := i \rightarrow o$
 $NP := i$
 $S := o$
 $J := j$
 $U := \lambda x. \mathbf{unicorn} x$
 $A := \lambda p. \lambda q. \exists x. p x \wedge q x$
 $S := \lambda p. \lambda x. \mathbf{try} x (\lambda y. p (\lambda z. \mathbf{find} y z))$

We have that:

$$\mathcal{L}_1(\text{S (A U) J}) = \text{/John/} + \text{/seeks/} + \text{/a/} + \text{/unicorn/}$$

$$\mathcal{L}_2(\text{S (A U) J}) = \text{try j}(\lambda x. \exists y. \text{unicorn } y \wedge \text{find } x y)$$

$$\mathcal{L}_1(\text{A U}(\lambda x. \text{S}(\lambda k. k x) \text{J})) = \text{/John/} + \text{/seeks/} + \text{/a/} + \text{/unicorn/}$$

$$\mathcal{L}_2(\text{A U}(\lambda x. \text{S}(\lambda k. k x) \text{J})) = \exists y. \text{unicorn } y \wedge \text{try j}(\lambda x. \text{find } x y)$$

A language-theoretic example

Abstract vocabulary:

$$\begin{aligned}
 A, L, S & : \text{ type} \\
 H & : (A \multimap A \multimap A \multimap S) \multimap S \\
 I & : L \multimap S \\
 E & : L \\
 C & : A \multimap L \multimap L
 \end{aligned}$$

Lexicon:

$$\begin{aligned}
 A, L, S & := \text{ string} \\
 H & := \lambda f. f /a/ /b/ /c/ \\
 I & := \lambda f. \lambda x. f x \\
 E & := \epsilon \\
 C & := \lambda x. \lambda y. x + y
 \end{aligned}$$

Typically:

$$H (\lambda x_{11} x_{12} x_{13}. H (\lambda x_{21} x_{22} x_{23}. \dots I (C x_{ij} (C x_{kl} \dots (C x_{mn} E) \dots))) \dots)) : S$$

Some Key Properties

Curry-Howard isomorphism

Coherence theorem

Principal typing

Subject reduction

Subject expansion

Constructing a Parsing Algorithm

Back to the example

$$\begin{array}{ll}
 H & := \lambda f. f (\lambda z. a z) (\lambda z. b z) (\lambda z. c z) : (A \multimap A \multimap A \multimap S) \multimap S \\
 I & := \lambda f. \lambda x. f x : L \multimap S \\
 E & := \lambda x. x : L \\
 C & := \lambda x. \lambda y. \lambda z. x (y z) : A \multimap L \multimap L
 \end{array}$$

$$A, L, S := s \multimap s$$

$\lambda z. a (c (b (a (b (c z)))))) ?$

A first non deterministic algorithm

1. Try to prove S using the types of the abstract constants as proper axioms.

I.e, prove S using $(A \multimap A \multimap A \multimap S) \multimap S$, $L \multimap S$, L , and $A \multimap L \multimap L$.

2. By the Curry-Howard isomorphism, you have constructed a term of the abstract language. Apply the lexicon to this term.

3. Check whether the resulting object term is equal to the term you have to parse.

The Coherence Theorem comes in

1. Specialize the object signature by distinguishing between the different occurrences of a same object constant in the term to be parsed:

$$a_1 : s_5 \multimap s_6$$

$$a_2 : s_2 \multimap s_3$$

$$b_1 : s_3 \multimap s_4$$

$$b_2 : s_1 \multimap s_2$$

$$c_1 : s_4 \multimap s_5$$

$$c_2 : s_0 \multimap s_1$$

$$\lambda z. a_1 (c_1 (b_1 (a_2 (b_2 (c_2 z)))))) : s_0 \multimap s_6$$

2. Specialize the lexical entries accordingly:

$$\lambda f. f (\lambda z. a_1 z) (\lambda z. b_1 z) (\lambda z. c_1 z) : \dots$$

$$\lambda f. f (\lambda z. a_1 z) (\lambda z. b_1 z) (\lambda z. c_2 z) : \dots$$

$$\dots : \dots$$

3. Try to prove $\langle S, s_0 \multimap s_6 \rangle$ using:

$$\begin{aligned}
 & \langle (A \multimap A \multimap A \multimap S) \multimap S, \\
 & \quad ((s_5 \multimap s_6) \multimap (s_3 \multimap s_4) \multimap (s_4 \multimap s_5) \multimap (s_0 \multimap s_0)) \multimap (s_0 \multimap s_0) \rangle \\
 & \langle (A \multimap A \multimap A \multimap S) \multimap S, \\
 & \quad ((s_5 \multimap s_6) \multimap (s_3 \multimap s_4) \multimap (s_4 \multimap s_5) \multimap (s_0 \multimap s_1)) \multimap (s_0 \multimap s_1) \rangle \\
 & \quad \vdots \\
 & \langle (A \multimap A \multimap A \multimap S) \multimap S, \\
 & \quad ((s_5 \multimap s_6) \multimap (s_3 \multimap s_4) \multimap (s_0 \multimap s_1) \multimap (s_0 \multimap s_0)) \multimap (s_0 \multimap s_0) \rangle \\
 & \langle (A \multimap A \multimap A \multimap S) \multimap S, \\
 & \quad ((s_5 \multimap s_6) \multimap (s_3 \multimap s_4) \multimap (s_0 \multimap s_1) \multimap (s_0 \multimap s_1)) \multimap (s_0 \multimap s_1) \rangle \\
 & \quad \vdots \\
 & \langle L \multimap S, (s_0 \multimap s_0) \multimap (s_0 \multimap s_0) \rangle \\
 & \langle L \multimap S, (s_0 \multimap s_1) \multimap (s_0 \multimap s_1) \rangle \\
 & \quad \vdots
 \end{aligned}$$

Eliminating redundancies

Consider the following pair:

$$\langle (A \multimap A \multimap A \multimap S) \multimap S, ((s_5 \multimap s_6) \multimap (s_3 \multimap s_4) \multimap (s_4 \multimap s_5) \multimap (s_0 \multimap s_0)) \multimap (s_0 \multimap s_0) \rangle$$

The shape of the specialized object type is completely specified by the grammar. The only relevant information is given by the indices.

Replace the above pair by the following formula:

$$(A[5, 6] \multimap A[3, 4] \multimap A[4, 5] \multimap S[0, 0]) \multimap S[0, 0]$$

Principal typing

Factorize the several formulas coming from a given lexical entry,

$$\begin{aligned}
 & (A[5, 6] \multimap A[3, 4] \multimap A[4, 5] \multimap S[0, 0]) \multimap S[0, 0] \\
 & (A[5, 6] \multimap A[3, 4] \multimap A[4, 5] \multimap S[0, 1]) \multimap S[0, 1] \\
 & \quad \vdots \\
 & (A[5, 6] \multimap A[3, 4] \multimap A[0, 1] \multimap S[0, 0]) \multimap S[0, 0] \\
 & (A[5, 6] \multimap A[3, 4] \multimap A[0, 1] \multimap S[0, 1]) \multimap S[0, 1] \\
 & \quad \vdots
 \end{aligned}$$

as follows:

$$a[i, j], b[k, l], c[m, n] \vdash (A[i, j] \multimap A[k, l] \multimap A[m, n] \multimap S[o, p]) \multimap S[o, p]$$

We end up with the following proof search problem:

Formulas coming from the lexicon:

$$\begin{aligned} & a[i, j], b[k, l], c[m, n] \vdash (A[i, j] \multimap A[k, l] \multimap A[m, n] \multimap S[o, p]) \multimap S[o, p] \\ & \vdash L[i, j] \multimap S[i, j] \\ & \vdash L[i, i] \\ & \vdash A[i, j] \multimap L[k, i] \multimap L[k, j] \end{aligned}$$

Query (coming from the term to be parsed):

$$a[5, 6], c[4, 5], b[3, 4], a[2, 3], b[1, 2], c[0, 1] \vdash S[0, 6]$$

Correctness and Completeness

Correctness : by subject reduction.

Completeness : by subject expansion.

The second-order case

— Kanazawa's original construction

- Allow the lexicons to be compiled into datalog programs.
- Polynomiality of 2nd-order ACGs.
- Optimizations techniques are known.
- CFG, TAG, LCFRS, ... as 2nd-order ACGs.