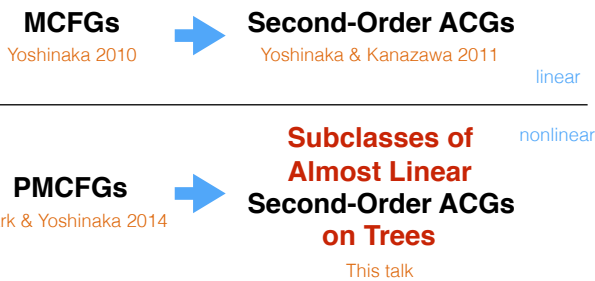


Distributional Learning and Context/Substructure Enumerability in Nonlinear Tree Grammars

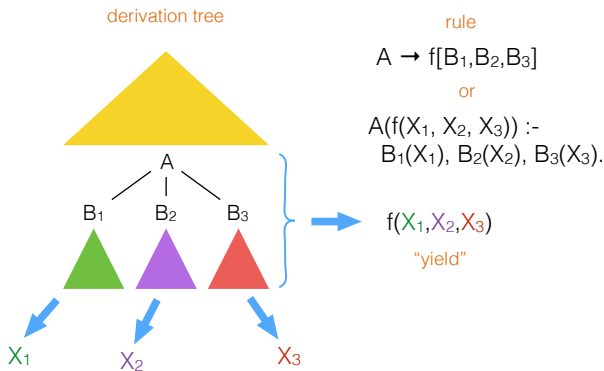
Makoto Kanazawa, National Institute of Informatics & SOKENDAI
Ryo Yoshinaka, Kyoto University

Generalizing Distributional Learning



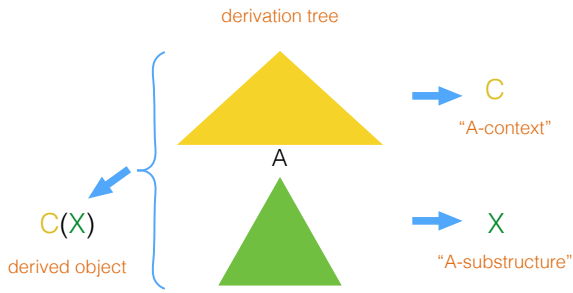
Almost linear second-order ACGs with bounded degree of nonlinearity either on the context side or the substructure side.
Will use the next 16 slides to give the background.

“Context-Free” Grammar Formalisms



4

Contexts and Substructures



Extraction of a subtree from a derivation tree induces a decomposition of the derived object into a context and a substructure.

5

Possible Contexts and Substructures

\mathcal{G} : class of grammars

\mathcal{C} : set of possible contexts for grammars in \mathcal{G}

\mathcal{S} : set of possible substructures for grammars in \mathcal{G}

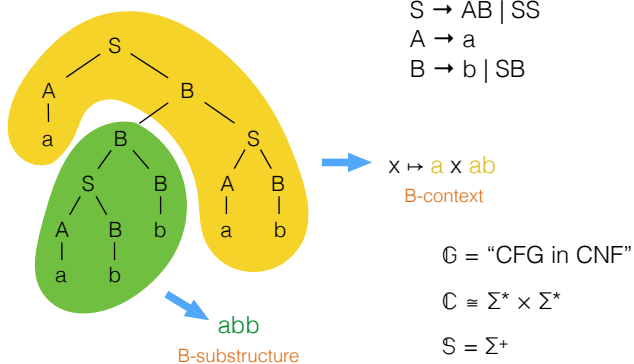
Each grammar G in \mathcal{G} determines a relation:

$$\{ (C, X) \in \mathcal{C} \times \mathcal{S} \mid C(X) \in L(G) \}$$

Note: $C(X) \in L(G)$ does not imply G has a nonterminal A such that C is an A -context and X is an A -substructure.

6

Example: CFG in CNF



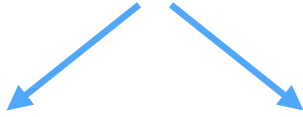
A context is a function that takes a string and wraps two strings around it.

Distributional Learning

7

“Tractable” means certain parameters of grammars are fixed.

Hypothesis space:
“tractable” grammar class \mathcal{G}



Dual approach:
target grammars in \mathcal{G} with
m-finite context property

Primal approach:
target grammars in \mathcal{G} with
m-finite kernel property

Finite Context Property

8

Let’s look at the dual approach, the simpler of the two approaches.

- G has the **m-FCP** \Leftrightarrow for each nonterminal A of G , there is a finite set \mathbf{C}_A of A -contexts such that
 - for all $X \in S$,
 $C(X) \in L(G)$ for all $C \in \mathbf{C}_A \Rightarrow X$ is an A -substructure
 - $|\mathbf{C}_A| \leq m$
- As a consequence, for each rule $A \rightarrow f[B_1, \dots, B_n]$ of G and for all $X_1, \dots, X_n \in S$,

$$\bigwedge_{1 \leq i \leq n} (C(X_i) \in L(G) \text{ for all } C \in \mathbf{C}_{B_i}) \Rightarrow C(f(X_1, \dots, X_n)) \in L(G) \text{ for all } C \in \mathbf{C}_A$$

CFG with 1-FCP

9

Each rule is valid under the given assignment of context sets to nonterminals.
(Removing some rules breaks 1-FCP, but validity of each rule is intact as long as the language stays the same.)

$S \rightarrow AB \mid SS$ $\mathbf{C}_S = \{ (\varepsilon, \varepsilon) \}$
 $A \rightarrow a \mid AS \mid SA$ $\mathbf{C}_A = \{ (\varepsilon, b) \}$
 $B \rightarrow b \mid SB \mid BS$ $\mathbf{C}_B = \{ (a, \varepsilon) \}$

x is an S -substructure $\Leftrightarrow nf(x) = \varepsilon$ $ab \succ \varepsilon$
 x is an A -substructure $\Leftrightarrow nf(x) = a$
 x is a B -substructure $\Leftrightarrow nf(x) = b$

10

Distributional Learning Algorithm (Dual Approach)

D: finite positive data (received up to some point)

$C|_D = \{ C \in \mathcal{C} \mid C(X) \in D \text{ for some } X \in \mathcal{S} \}$ observed contexts

Construct rules

$$\mathbf{C}_0 \rightarrow f[\mathbf{C}_1, \dots, \mathbf{C}_n]$$

such that

- \mathbf{C}_i is a nonempty finite subset of $C|_D$ with $|\mathbf{C}_i| \leq m$
- f is a function “observed” in D

D is the positive data received up to the last moment when the hypothesis grammar was found to undergenerate.

11

Distributional Learning Algorithm (Dual Approach)

$E \supseteq D$: all positive data received so far

$\mathcal{S}|_E = \{ X \in \mathcal{S} \mid C(X) \in E \text{ for some } C \in \mathcal{C} \}$ observed substructures

Hypothesize grammar using those rules

$$\mathbf{C}_0 \rightarrow f[\mathbf{C}_1, \dots, \mathbf{C}_n]$$

constructed from D that are **valid for** $\mathcal{S}|_E$:

for all $X_1, \dots, X_n \in \mathcal{S}|_E$,

$$\bigwedge_{1 \leq i \leq n} (C(X_i) \in L(G) \text{ for all } C \in \mathbf{C}_i) \Rightarrow$$

$$C(f(X_1, \dots, X_n)) \in L(G) \text{ for all } C \in \mathbf{C}_0$$

use membership oracle

Membership queries are restricted to combinations of bits of observed positive data. (Not quite?)

12

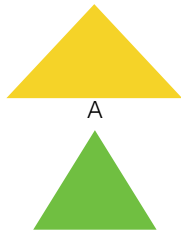
Context/Substructure Enumerability

- A grammar class \mathcal{G} is
 - **context-enumerable** if $C|_D$ can be enumerated in polynomial time in the size of D
 - **substructure-enumerable** if $\mathcal{S}|_D$ can be enumerated in polynomial time in the size of D
- Both are needed to achieve polynomial update time of distributional learning.

Clearly satisfied by $\mathcal{G} = \text{“CFG in CNF”}$.

MCFGs (of bounded dimension and rank)

derivation tree



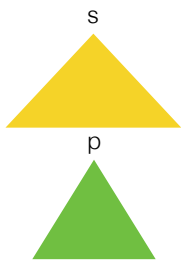
$(x_1, \dots, x_r) \mapsto v_0 x_1 v_1 \dots x_r v_r$
 regular r -variable pattern \equiv
 $(r+1)$ -tuple of strings

$(w_1, \dots, w_r) \in (\Sigma^*)^r$
 r -tuple of strings

context-enumerable	✓
substructure-enumerable	✓

Second-Order (Linear) ACGs

derivation tree



N
 closed linear λ -term
 of type $\mathcal{H}(p) \rightarrow \mathcal{H}(s)$

M
 closed linear λ -term of type $\mathcal{H}(p)$

context-enumerable	✓
substructure-enumerable	✓

With a similar bound on certain parameters (possible object realizations of abstract types of constants limited to a fixed finite set). The type $\mathcal{H}(p)$ is taken from a fixed finite set.

k-copying PMCFGs

$S(x_1 x_2 a) :- A(x_1, x_2)$
 $A(x_1 x_2 a, x_2 a a) :- A(x_1, x_2)$ 2-copying
 $A(\epsilon, \epsilon) :-$

$(x_1, x_2) \mapsto$
 $x_1 x_2 a x_2 a a x_2 a a a a a$
 3-copying context



k -copying means each rule can make at most k copies of a single string. A context need not be k -copying.

k-copying PMCFGs (of bounded dimension and rank)

derivation tree



A



→ $(x_1, \dots, x_r) \mapsto v_0 x_{i_1} v_1 \dots x_{i_n} v_n$
non-regular r-variable pattern
(with no bound on n)

→ $(w_1, \dots, w_r) \in (\Sigma^*)^r$
r-tuple of strings

context-enumerable	✗
substructure-enumerable	✓

There is no bound on the number of copies a context makes of a component of a substructure.

k-copying IO CFTGs

derivation tree



A



→ $t[x_1, \dots, x_r] \mapsto u[t[v_1, \dots, v_r]]$
one-variable tree pattern +
r-tuple of trees

→ $t[x_1, \dots, x_r] \in T_\Delta(X_r)$
r-variable tree pattern

context-enumerable	✗
substructure-enumerable	✗

The dual approach.

Narrowing the Learning Target for k-copying PMCFGs (C&Y 2014)

\mathcal{G} = the class of k-copying PMCFGs bounded dimension, rank

$\mathcal{C}_k = \{ C \in \mathcal{C} \mid C \text{ is } k\text{-copying} \}$

$\mathcal{C}_k|_D = \{ C \in \mathcal{C}_k \mid C(X) \in D \text{ for some } X \in \mathcal{S} \}$ poly-time enumerable

Use finite subsets \mathbf{C} of $\mathcal{C}_k|_D$ (with $|\mathbf{C}| \leq m$) as nonterminals

Validate candidate rules using $\mathcal{S}|_E$ (where $D \subseteq E$)

Works for k-copying PMCFGs with **(k,m)-FCP**:

- each nonterminal A is characterized by a finite set \mathbf{C}_A of **k-copying** A-contexts (with $|\mathbf{C}_A| \leq m$)

Goal of This Work

19

Use almost linear second-order ACGs on trees.

- Define large classes of nonlinear grammars that are context- or substructure-enumerable.
- Show that PMCFG-type distributional learning works for these classes.

Almost Linear Second-Order ACGs on Trees

20

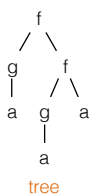
Generalizes IO CFTGs.

- Use **almost linear λ -terms** as functions in rules.
- A general formalism encompassing all known efficiently recognizable “context-free” formalisms on trees.
- Equivalent to
 - attribute grammars outputting trees
 - MSO-definable tree-to-term-graph transductions + unfolding

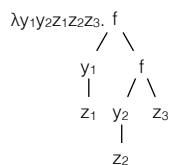
Almost Linear λ -terms

21

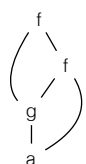
A nonlinear λ -term sometimes can be made to look linear by sharing of subterms (of atomic type). In that case, the λ -term is said to be almost linear.



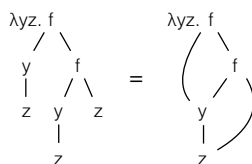
tree



linear λ -term



tree with sharing



almost linear λ -term

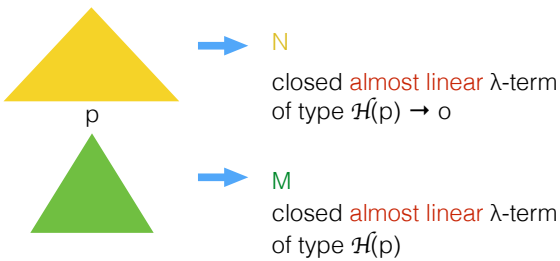
Almost Linear ACGs

$s(X(\lambda z^o.fzz)) \text{ :- } p(X)$ $f : o \rightarrow o \rightarrow o$
 $p(\lambda y^{o \rightarrow o}.X(\lambda z^o.y(fzz))) \text{ :- } p(X)$ $a : o$
 $p(\lambda y^{o \rightarrow o}.ya) \text{ :-}$ $\mathcal{H}(s) = o$
 $\mathcal{H}(p) = (o \rightarrow o) \rightarrow o$

s $f(f(faa)(faa))(f(faa)(faa))$ perfect binary trees
 $|$
 p $\lambda y^{o \rightarrow o}.y(f(faa)(faa))$
 $|$
 p $\lambda y^{o \rightarrow o}.y(faa)$
 $|$
 p $\lambda y^{o \rightarrow o}.ya$

Almost Linear Second-Order ACGs

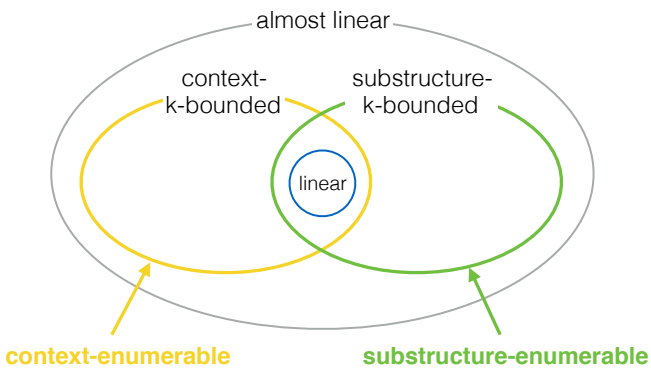
derivation tree



context-enumerable	✗
substructure-enumerable	✗

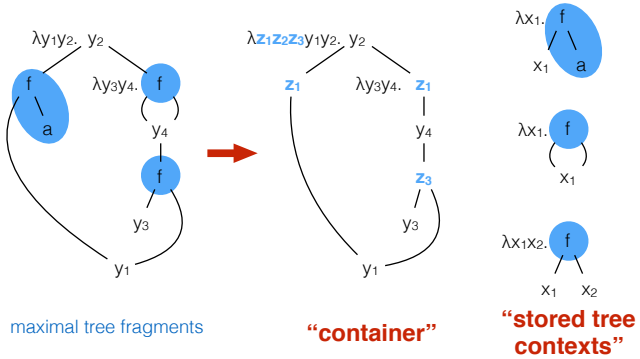
Neither context- nor substructure-enumerable even with a bound on certain parameters.
Extends IO CFTGs.

Almost Linear ACGs with Bounded Nonlinearity



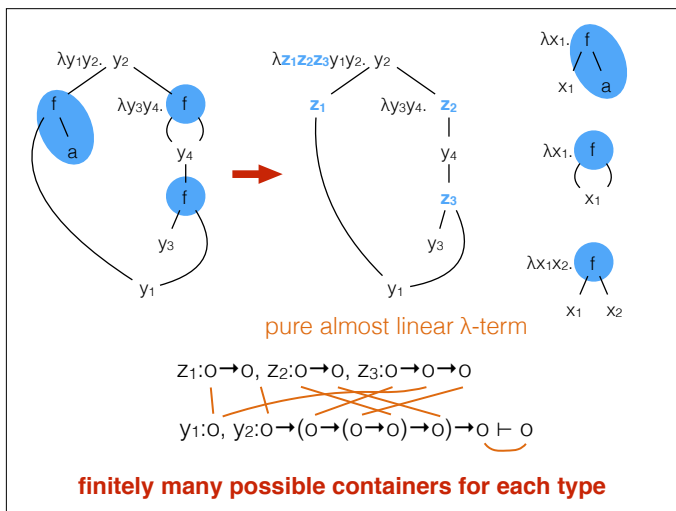
Limit the “amount” of nonlinearity that shows up in substructures or in contexts.

Decomposition of Almost Linear λ -terms over a Tree Signature



Preparation for the definition of k -boundedness of an almost linear λ -term.

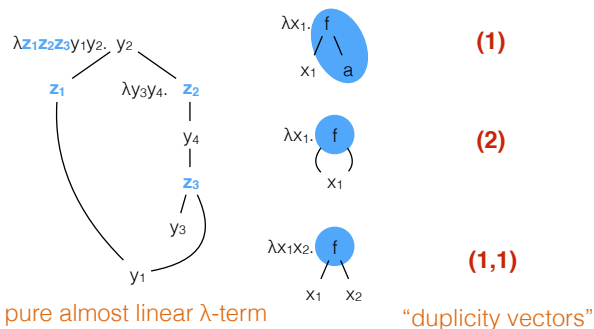
Extract maximal tree fragments from a given almost linear λ -term in such a way that the remaining pure λ -term is almost linear.



The number of stored tree contexts is bounded by the number of positive atomic type occurrences in the type of the input λ -term.

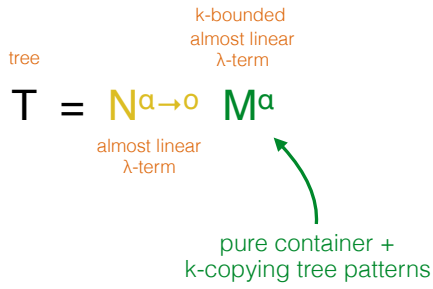
The arity of each stored tree context is bounded by the number of negative atomic type occurrences in the type of the input λ -term.

Degree of Nonlinearity



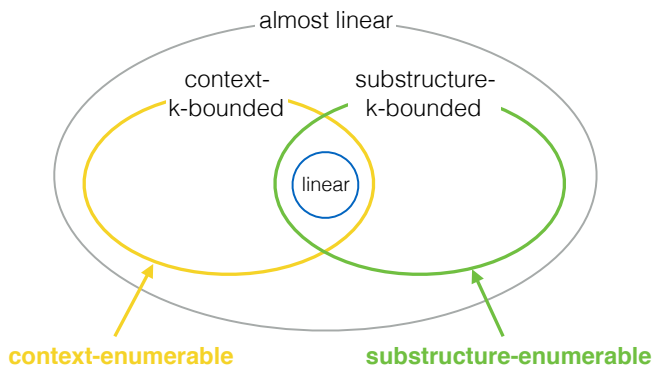
M is k -bounded $\Leftrightarrow M$'s duplicity vectors $\in \{1, \dots, k\}^*$

Extracting k-bounded λ-terms from a Tree



Such M are polynomial-time enumerable.

Almost Linear ACGs with Bounded Nonlinearity



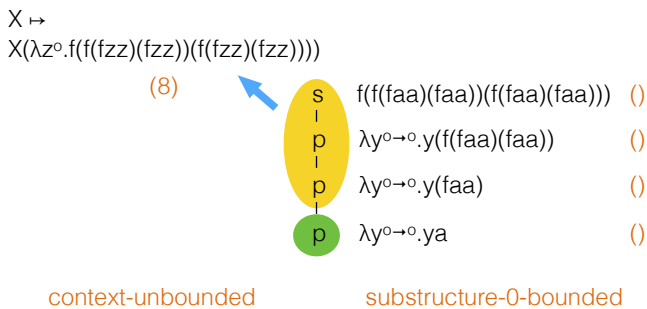
Substructure-k-bounded: Every substructure the grammar provides as the yield of a subtree of a derivation tree is a k-bounded almost linear λ-term.

Context-k-bounded: Every context the grammar provides as the yield of a derivation tree with a hole is a k-bounded almost linear λ-term.

$s(X(\lambda z^o.fzz)) :- p(X) \quad (2)$
 $p(\lambda y^{o \rightarrow o}.X(\lambda z^o.y(fzz))) :- p(X) \quad (2)$
 $p(\lambda y^{o \rightarrow o}.ya) :- \quad ()$

rule-2-bounded

An example of a substructure-k-bounded almost linear ACG.



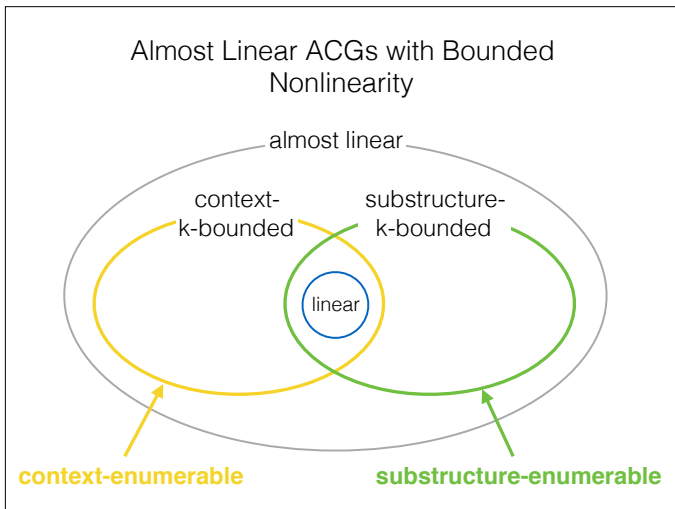
31

An example that is unbounded on both sides.

$s(Xa) :- p(X)$ $p(\lambda z^0.f(Xz)(Xz)) :- p(X)$ $p(\lambda z^0.z) :-$	$()$ (2) $-$												
rule-2-bounded													
$X \mapsto$ $f(f(Xa)(Xa))(f(Xa)(Xa))$ $(4), ()$	<table style="margin-left: auto; margin-right: auto;"> <tr> <td style="text-align: center; padding: 2px;">s</td> <td style="padding: 2px;">f(faa)(faa)</td> <td style="padding: 2px;">()</td> </tr> <tr> <td style="text-align: center; padding: 2px;">p</td> <td style="padding: 2px;">$\lambda z^0.f(fzz)(fzz)$</td> <td style="padding: 2px;">(4)</td> </tr> <tr> <td style="text-align: center; padding: 2px;">p</td> <td style="padding: 2px;">$\lambda z^0.fzz$</td> <td style="padding: 2px;">(2)</td> </tr> <tr> <td style="text-align: center; padding: 2px;">p</td> <td style="padding: 2px;">$\lambda z^0.z$</td> <td style="padding: 2px;">-</td> </tr> </table>	s	f(faa)(faa)	()	p	$\lambda z^0.f(fzz)(fzz)$	(4)	p	$\lambda z^0.fzz$	(2)	p	$\lambda z^0.z$	-
s	f(faa)(faa)	()											
p	$\lambda z^0.f(fzz)(fzz)$	(4)											
p	$\lambda z^0.fzz$	(2)											
p	$\lambda z^0.z$	-											
context-unbounded	substructure-unbounded												

32

Are these classes decidable?



33

k-threshold Profile of an Almost Linear λ -term

$$\text{prof}(M) = \left(\begin{array}{c} \text{container} \\ M^0, \quad w_1, \dots, w_l \end{array} \right)$$

duplicity vectors
 $\in (\mathbb{N} - \{0\})^*$

$$\text{prof}_k(M) = \left(\begin{array}{c} \text{container} \\ M^0, \quad w_1, \dots, w_l \end{array} \right)$$

duplicity vectors
 $\in \{1, \dots, k, \infty\}^*$

Lemma. For each type α , there are only finitely many k-threshold profiles of closed almost linear λ -terms (over a tree signature) of type α .

Deciding k-boundedness

Lemma. If $\text{prof}_k(M_i) = \text{prof}_k(M_i')$ for $i = 1, \dots, n$, then $\text{prof}_k(NM_1 \dots M_n) = \text{prof}_k(NM_1' \dots M_n')$.

Lemma. The derivation trees of an almost linear ACG that involve only k-bounded substructures form a regular set.

Lemma. The derivation trees of an almost linear ACG that involve only k-bounded contexts form a regular set.

The class of substructure-k-bounded almost linear ACGs is a decidable class.

So is the class of context-k-bounded almost linear ACGs.

Can take either of these classes (with certain parameters fixed) as hypothesis space of distributional learning.

PMCFG-Type Dual Distributional Learning of Substructure-k-bounded Almost Linear ACGs

\mathbb{G} = the class of almost linear ACGs that are rule- and substructure-k-bounded with certain parameters fixed

$\mathbb{C}_k = \{ C \in \mathbb{C} \mid C \text{ is } k\text{-bounded} \}$

$\mathbb{C}_{k|D} = \{ C \in \mathbb{C}_k \mid C(X) \in D \text{ for some } X \in \mathcal{S} \}$ poly-time enumerable

Use finite subsets \mathbf{C} of $\mathbb{C}_{k|D}$ (with $|\mathbf{C}| \leq m$) as “nonterminals”

Validate candidate rules using $\mathcal{S}|_E$ (where $D \subseteq E$)

Target grammars in \mathbb{G} with **(k,m)-FCP**:

- each “nonterminal” p is characterized by a finite set \mathbf{C}_p of **k-bounded** λ -terms of type $\mathcal{H}(p) \rightarrow o$ (with $|\mathbf{C}_p| \leq m$)

There is no guarantee that the hypothesized grammar is substructure-k-bounded.

This can cause overgeneration.

Ensuring Substructure-k-boundedness

- Use (\mathbf{C}, π) as nonterminals where
 - \mathbf{C} is a finite subset of $\mathbb{C}_{k|D}$ consisting of λ -terms of type α (with $|\mathbf{C}| \leq m$)
 - π is a k-threshold profile for type α
- Construct rules

$(\mathbf{C}_0, \pi_0)(P^{\alpha_1 \rightarrow \dots \rightarrow \alpha_n \rightarrow \alpha_0} X_1 \dots X_n) :- (\mathbf{C}_1, \pi_1)(X_1), \dots, (\mathbf{C}_n, \pi_n)(X_n)$

such that

 - P is a k-bounded almost linear λ -term “observed” in D
 - $\pi_0 = P\pi_1 \dots \pi_n$.
- Validate for $X_i \in \mathcal{S}|_E \cap \text{prof}_k^{-1}(\pi_i)$.

The second bullet point guarantees substructure-k-boundedness.

Profile-Sensitive (k,m)-FCP

37

- Each “nonterminal” p of G is characterized by finite sets $\mathbf{C}_{p,\pi}$ of \mathbb{C}_k (with $|\mathbf{C}_{p,\pi}| \leq m$), one for each profile π for $\mathcal{H}(p)$:

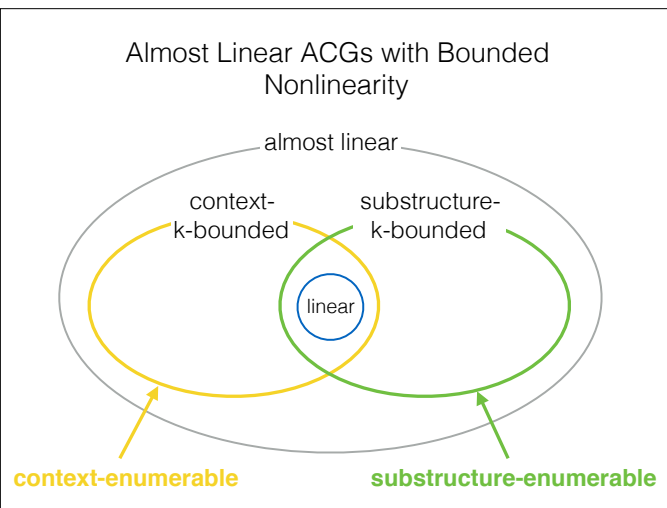
for all $X \in \mathcal{S} \cap \text{prof}_k^{-1}(\pi)$,

$C(X) \in L(G)$ for all $C \in \mathbf{C}_{p,\pi} \Rightarrow X$ is a p -substructure

This is a more natural notion of FCP corresponding to the validation in the learning algorithm.

Almost Linear ACGs with Bounded Nonlinearity

38



“PMCFG-type” distributional learning is possible for both these classes (more complicated for the context-k-bounded grammars). Leads to a “profile-sensitive” notion of FCP (and of FKP).

39

