

Apprentissage de langages non réguliers

Soutenance de stage

Jean-Marie Madiot

LIF de Marseille, ENS Lyon

9 septembre 2009

1 Introduction

2 CBFG

- Définitions
- Différence avec les grammaires hors-contexte
- CBFG exacte
- Exemple
- Inférence, algorithme

3 Expérimentations

- Protocole expérimental
- Choix de la classe de test
- Génération des données

4 Résultats

5 Conclusion

6 Questions ?

Apprentissage de langages naturels

Langages
learning

Jean-Marie
Madiot

Introduction

CBFG

Definition

Comparison

Exact CBFG

Example

Inference

Experiments

Protocol

Testing set

Generation

Results

Conclusion

Questions

- Heuristiques
Peu de formalisme (garantie)
- Classes existantes apprenables (réguliers, very simple)
Trop petites
- Classes puissantes (context-free, mildly context-sensitive)
Difficilement apprenables

Approche

Relier apprenabilité et formalisme

Langages
learning

Jean-Marie
Madiot

Introduction

CBFG

Definition
Comparison
Exact CBFG
Example
Inference

Experiments

Protocol
Testing set
Generation

Results

Conclusion

Questions

Apprentissage : basé sur l'observable

- distribution
- contextes
- ...

Idée : formaliser à l'aide des contextes

CBFG

Contextual Binary Feature Grammar

Langages
learning

Jean-Marie
Madiot

Introduction

CBFG

Definition

Comparison

Exact CBFG

Example

Inference

Experiments

Protocol

Testing set

Generation

Results

Conclusion

Questions

une CBFG G est la donnée de :

- Σ , alphabet
- F , ensemble fini de contextes
- P , ensemble de règles $A \rightarrow B \cdot C$ ($A, B, C \subseteq F$)
- P_L , ensemble de règles $A \rightarrow a$ ($A \subseteq F, a \in \Sigma$)

Principe : $A \rightarrow B \cdot C$ signifie : si u, v est **associé** à tous les contextes de B, C , on associe à uv tous les contextes de A .

« on associe (a, b) à u » : $(a, b) \in f_G(u)$.

« u est dans le langage » : $(\lambda, \lambda) \in f_G(u)$

CBFG

Contextual Binary Feature Grammar

Langages
learning

Jean-Marie
Madiot

Introduction

CBFG

Definition

Comparison

Exact CBFG

Example

Inference

Experiments

Protocol

Testing set

Generation

Results

Conclusion

Questions

$$f_G(\lambda) = \emptyset \quad (1)$$

$$f_G(a) = \bigcup_{x \rightarrow a} x \quad \text{si } |a| = 1 \quad (2)$$

$$f_G(w) = \bigcup_{u,v:uv=w} \bigcup_{\substack{x \rightarrow yz: \\ y \subseteq f_G(u) \wedge \\ z \subseteq f_G(v)}} x \quad \text{si } |w| > 1. \quad (3)$$

Analyse ascendante

Parsing

Langages
learning

Jean-Marie
Madiot

Introduction

CBFG

Definition

Comparison

Exact CBFG

Example

Inference

Experiments

Protocol

Testing set

Generation

Results

Conclusion

Questions

Calcul de f_G sur les facteurs de w .

- $|w| = 1, w = a :$

$$f_G(a) = \bigcup_{x \rightarrow a} x$$

- $|w|$ quelconque : pour tous $u, v, w = uv$, on inclut les x tels que $x \rightarrow yz$, et tels que $f_G(u)$ contienne y (et $f_G(v)$ z)

Différence avec les grammaires hors-contexte

Langages
learning

Jean-Marie
Madiot

Introduction

CBFG

Definition

Comparison

Exact CBFG

Example

Inference

Experiments

Protocol

Testing set

Generation

Results

Conclusion

Questions

Hors-contextes : génératifs

CBFG : contraintes

$A, C, D \subseteq F$, ensemble de contextes

- $A (a^* b^n c^n)$
- $C (a^n b^n c^*)$
- $D (d)$

$$S \rightarrow A \cup C \cdot D$$

correspond à $a^n b^n c^n d$. (contextuel)

CBFG exacte

Langages
learning

Jean-Marie
Madiot

Introduction

CBFG

Definition

Comparison

Exact CBFG

Example

Inference

Experiments

Protocol

Testing set

Generation

Results

Conclusion

Questions

F est abstrait :

les contextes de F n'ont pas de sens par rapport au langage
la **structure** du langage ne correspond pas à son **apparence**

Intérêt d'apprentissage : quand les contextes de F coïncident
avec ceux du langage

$(a, \lambda) \in F$ peut être un contexte de b , sans que $f_G(b)$
contienne (a, λ) .

CBFG exacte : quand $f_G(u)$ contient les contextes de F qui
apparaissent autour de u dans le langage

Un langage qui n'admet pas de CBFG exacte

Langages
learning

Jean-Marie
Madiot

Introduction

CBFG

Definition

Comparison

Exact CBFG

Example

Inference

Experiments

Protocol

Testing set

Generation

Results

Conclusion

Questions

$$L = \{a^n b^m, n > m > 0\} \cup \{a^n c, n > 0\}$$

Si L admet une CBFG exacte G :

$f_G(c) \subset f_G(b^k)$ car b^k a les contextes de c dans L .

$ac \in L$ donc il existe une règle :

$$\{\dots, (\lambda, \lambda), \dots\} \rightarrow D \cdot E$$

avec $D \subseteq f_G(a)$, $E \subseteq f_G(c) = f_G(b^k)$.

Donc $(\lambda, \lambda) \in f_G(ab^k)$.

Quelles classes de langages ?

Langages
learning

Jean-Marie
Madiot

Introduction

CBFG

Definition

Comparison

Exact CBFG

Example

Inference

Experiments

Protocol

Testing set

Generation

Results

Conclusion

Questions

La classe des CBFG :

- \supset hors-contextes
- \supset substituables
- = grammaires conjonctives

La classe apprenable (CBFG **exactes**) :

- \supset réguliers
- \supset substituables, *very simple*
- $\not\subseteq$ hors-contextes
- \subset grammaires conjonctives

Exemple : $a^n b^n$

Règles

Langages
learning

Jean-Marie
Madiot

Introduction

CBCFG

Definition

Comparison

Exact CBCFG

Example

Inference

Experiments

Protocol

Testing set

Generation

Results

Conclusion

Questions

$$F = \{(\lambda, \lambda), (a, \lambda), (\lambda, b), (aab, \lambda), (\lambda, abb)\}$$

- $\{(\lambda, b), (\lambda, abb)\} \rightarrow a$
- $\{(a, \lambda), (aab, \lambda)\} \rightarrow b$
- $\{(\lambda, \lambda)\} \rightarrow \{(\lambda, b)\}\{(aab, \lambda)\}$
- $\{(\lambda, \lambda)\} \rightarrow \{(\lambda, abb)\}\{(a, \lambda)\}$
- $\{(\lambda, b)\} \rightarrow \{(\lambda, abb)\}\{(\lambda, \lambda)\}$
- $\{(a, \lambda)\} \rightarrow \{(\lambda, \lambda)\}\{(\lambda, aab)\}$

Exemple : $a^n b^n$

Parsing

Languages
learning

Jean-Marie
Madiot

Introduction

CBFG

Definition

Comparison

Exact CBFG

Example

Inference

Experiments

Protocol

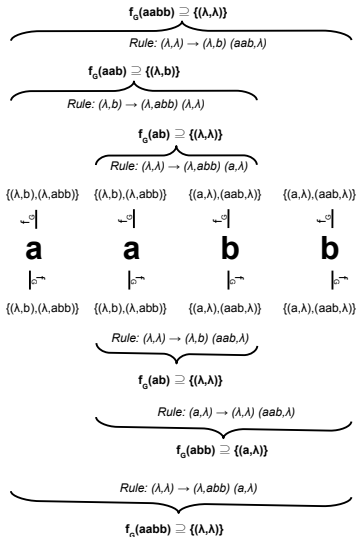
Testing set

Generation

Results

Conclusion

Questions



Construction d'une grammaire

Langages
learning

Jean-Marie
Madiot

Introduction

CBFG

Definition

Comparison

Exact CBFG

Example

Inference

Experiments

Protocol

Testing set

Generation

Results

Conclusion

Questions

K , ensemble de mots

F , ensemble de contextes

L , langage

$F_L(u)$: ensemble des contextes de F qui entourent u dans L

Construction d'une grammaire

Langages
learning

Jean-Marie
Madiot

Introduction

CBFG

Definition

Comparison

Exact CBFG

Example

Inference

Experiments

Protocol

Testing set

Generation

Results

Conclusion

Questions

K , ensemble de mots

F , ensemble de contextes

L , langage

P_L est l'ensemble des $F_L(a) \rightarrow a$

P est l'ensemble des $F_L(uv) \rightarrow F_L(u) \cdot F_L(v)$

on définit $G_0(K, L, F) = \langle F, P, P_L, \Sigma \rangle$

Construction d'une grammaire

Langages
learning

Jean-Marie
Madiot

Introduction

CBFG

Definition

Comparison

Exact CBFG

Example

Inference

Experiments

Protocol

Testing set

Generation

Results

Conclusion

Questions

$$G_0(K, L, F) = \langle F, P, P_L, \Sigma \rangle$$

Apprentissage d'un langage L donné :

F et de K donnent une grammaire $G_0(K, L, F)$.

Choix de F et K :

- $G_0(K, L, F)$ est croissant selon K
- $G_0(K, L, F)$ est décroissant selon F

Cet algorithme apprend les CBFG exactes.

Algorithme d'apprentissage

Langages
learning

Jean-Marie
Madiot

Introduction

CBFG

Definition

Comparison

Exact CBFG

Example

Inference

Experiments

Protocol

Testing set

Generation

Results

Conclusion

Questions

Utilisation d'un **oracle** : « Membership queries »

Principe :

S = les mots (positifs) en entrée

$K \subset$ les facteurs de S (*kernel*)

$F \subset$ tous les contextes de S

Combinaison de F et K et soumission à l'oracle ;

s'il existe un faux positif, on augmente F

s'il existe un faux négatif, on augmente K et F

Algorithme d'apprentissage

Langages
learning

Jean-Marie
Madiot

Introduction

CBFG

Definition

Comparison

Exact CBFG

Example

Inference

Experiments

Protocol

Testing set

Generation

Results

Conclusion

Questions

Algorithme 1 : Algorithme d'apprentissage de CBFG

Données : des mots, $\{w_1, w_2, \dots\}$, un oracle \mathcal{O}

Résultat : Une suite de CBFG G_1, G_2, \dots

$K \leftarrow \emptyset$ $F \leftarrow \{(\lambda, \lambda)\}$ $D \leftarrow \emptyset$ $G_0 = G_0(K, \mathcal{O}, F)$;

pour w_i **faire**

$D \leftarrow D \cup \{w_i\}$;

$S \leftarrow \text{Con}(D) \odot \text{Sub}(D)$;

si $\exists w \in S$ *tel que* $w \in L(G_{i-1}) \setminus L$ **alors**

$F \leftarrow \text{Con}(D)$;

fin

si $\exists w \in S$ *tel que* $w \in L \setminus L(G_{i-1})$ **alors**

$K \leftarrow \text{Sub}(D)$;

$F \leftarrow \text{Con}(D)$;

fin

 Renvoyer $G_i = G_0(K, \mathcal{O}, F)$;

fin

Expérimentations : but

Langages
learning

Jean-Marie
Madiot

Introduction

CBFG

Definition

Comparison

Exact CBFG

Example

Inference

Experiments

Protocol

Testing set

Generation

Results

Conclusion

Questions

Convergence : prouvée
Algorithme : polynomial

Comportement ?

- temps raisonnable ?
- approche progressive/soudaine ?
- sur/sous-généralisation ?

Protocole expérimental

Langages
learning

Jean-Marie
Madiot

Introduction

CBFG

Definition
Comparison
Exact CBFG
Example
Inference

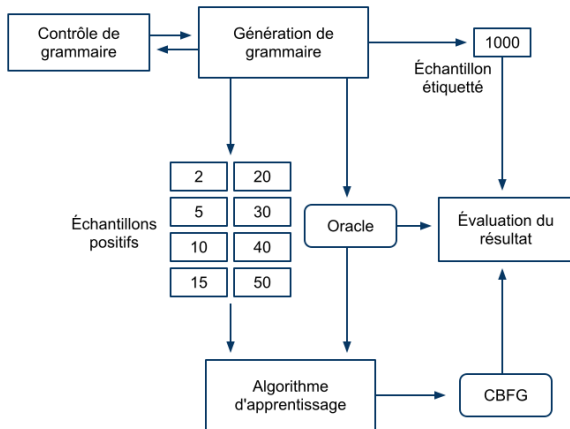
Experiments

Protocol
Testing set
Generation

Results

Conclusion

Questions



Choix de la classe de test

Langages
learning

Jean-Marie
Madiot

Introduction

CBFG

Definition

Comparison

Exact CBFG

Example

Inference

Experiments

Protocol

Testing set

Generation

Results

Conclusion

Questions

CBFG

- CBFG : non génératif, pas d'ESC
- exactitude indécidable

Mais les hors-contexte

- toujours un défi
- bonne variété parmi les CBFG
- connus
(comparaisons, résultats existants, notoriété)

Génération de grammaire

Hors contexte

Langages
learning

Jean-Marie
Madiot

Introduction

CBFG

Definition

Comparison

Exact CBFG

Example

Inference

Experiments

Protocol

Testing set

Generation

Results

Conclusion

Questions

- T : terminaux
- N : non-terminaux
- R : règles

$\#T$, $\#N$ et $\#R$ sont fixés.

Graphe de sommets T et N

Une règle : deux arêtes

On relie aléatoirement les sommets à l'aide d'arêtes.

On garde la composante connexe du graphe qui contient S .

Analyse de grammaire

Langages
learning

Jean-Marie
Madiot

Introduction

CBFG

Definition

Comparison

Exact CBFG

Example

Inference

Experiments

Protocol

Testing set

Generation

Results

Conclusion

Questions

Exploration du graphe à partir des terminaux

Calculable :

- L est fini (ou $L = \emptyset$)
→ génération d'une nouvelle grammaire
- des éléments sont inaccessibles
→ élimination des éléments

Incalculable :

- $L = \Sigma^*$
- L est régulier

Ensemble structurellement complet

Langages
learning

Jean-Marie
Madiot

Introduction

CBFG

Definition

Comparison

Exact CBFG

Example

Inference

Experiments

Protocol

Testing set

Generation

Results

Conclusion

Questions

Ensemble structurellement complet :

Ensemble de mots S tq :

$$\forall r \in R, \exists w \in L, w \text{ se dérive avec } r$$

Construction : exploration à partir des terminaux.

Idée de l'algorithme :

Calcul de l'ensemble des couples (règles accessibles, mots) pour chaque sommet

Fusions successives de ces ensembles

Détections des boucles (le langage est alors infini)

Génération de l'échantillon

Langages
learning

Jean-Marie
Madiot

Introduction

CBFG

Definition

Comparison

Exact CBFG

Example

Inference

Experiments

Protocol

Testing set

Generation

Results

Conclusion

Questions

Mots : *parsing* d'un mot aléatoire

- mots positifs : apprentissage
- mots étiquetés : test (minimum 40% de chaque)
- oracle : sous la forme d'une grammaire hors-contexte

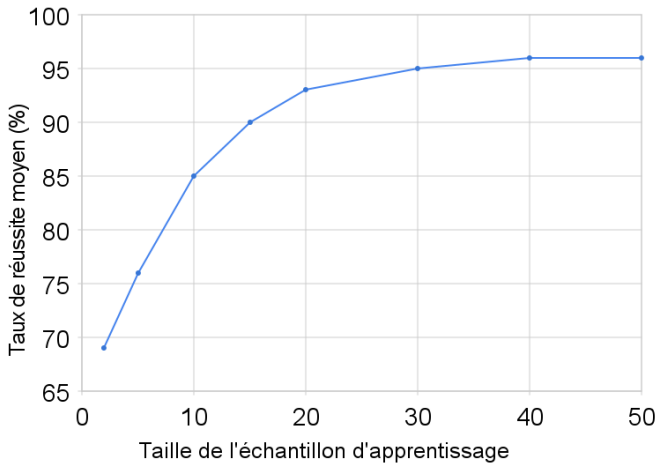
Tailles d'apprentissage : 2, 5, 10, 15, 20, 30, 40, 50

Même échantillon de test, de taille 1000

Passage à l'échelle : 10000

Résultats

Taux de réussite moyen



Langages
learning

Jean-Marie
Madiot

Introduction

CBFG

Definition
Comparison
Exact CBFG
Example
Inference

Experiments

Protocol
Testing set
Generation

Results

Conclusion

Questions

Résultats

Proportion au dessus d'un score donné

Langages
learning

Jean-Marie
Madiot

Introduction

CBFG

Definition

Comparison

Exact CBFG

Example

Inference

Experiments

Protocol

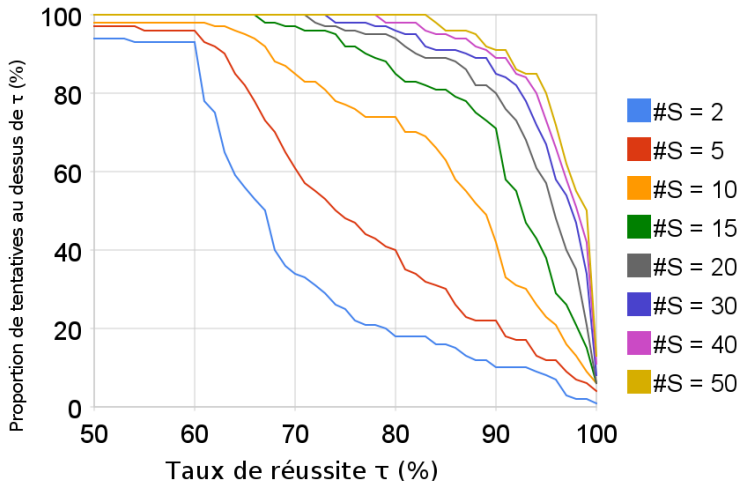
Testing set

Generation

Results

Conclusion

Questions



Résultats

Répartition des erreurs

Langages
learning

Jean-Marie
Madiot

Introduction

CBFG

Definition

Comparison

Exact CBFG

Example

Inference

Experiments

Protocol

Testing set

Generation

Results

Conclusion

Questions

#S	faux positifs	faux négatifs
2	2.3 %	28.1 %
5	3.3 %	20.0 %
10	2.4 %	12.1 %
15	2.6 %	6.5 %
20	1.6 %	4.6 %
30	1.1 %	3.6 %
40	0.69 %	3.0 %
50	0.55 %	2.5 %

Sous-généralisation (recherchée)

Conclusion

Langages
learning

Jean-Marie
Madiot

Introduction

CBFG

Definition

Comparison

Exact CBFG

Example

Inference

Experiments

Protocol

Testing set

Generation

Results

Conclusion

Questions

Problèmes

- ensemble structurellement complet
- nombre important de requêtes (naïf)

Justification de l'oracle

- langue maternelle
- « Zulu, interactive learning competition »

Questions ouvertes

- remplacer l'oracle
- préciser la classe des CBFG exactes
- sémantique d'une dérivation

Des questions ?

Langages
learning

Jean-Marie
Madiot

Introduction

CBFG

Definition

Comparison

Exact CBFG

Example

Inference

Experiments

Protocol

Testing set

Generation

Results

Conclusion

Questions