

Probabiliser les Grammaires de Propriétés

Philippe Blache & Stéphane Rauzy

Laboratoire Parole et Langage
CNRS & Aix-Marseille Université

- L'analyse en Grammaire de Propriétés
- Facteurs de complexité
- Quelles informations pour guider l'analyse ?
- Comment intégrer les probabilités ?

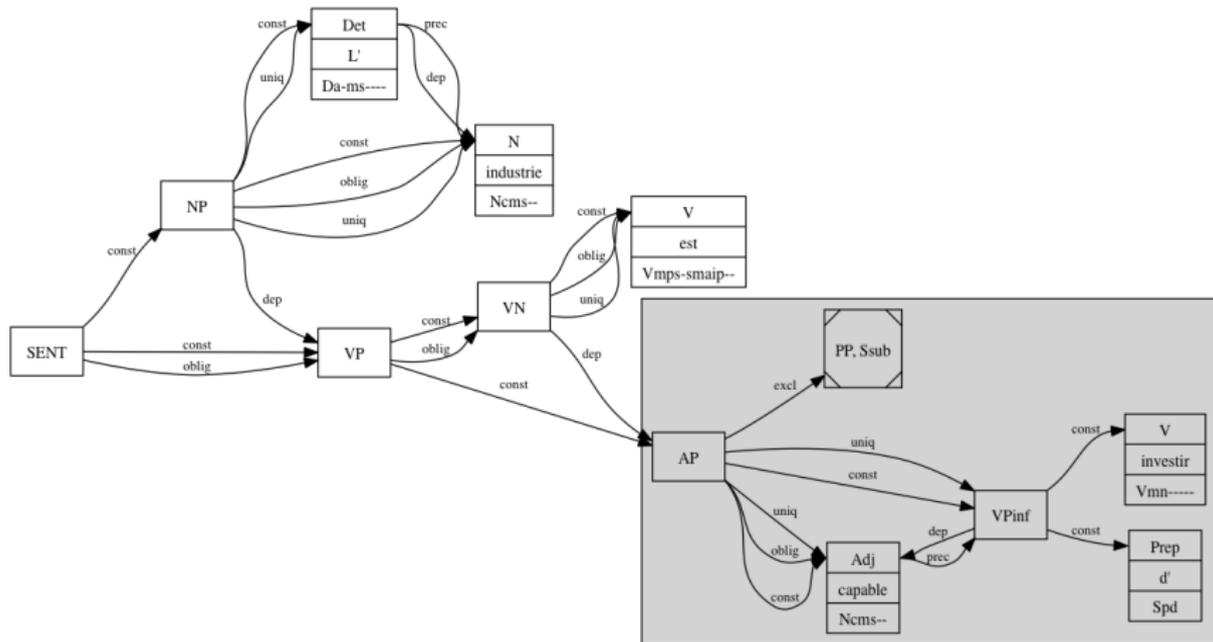
Les contraintes en GP

Obligation	$A : \Delta B$	au moins un B
Unicité	$A : B!$	au plus un B
Linéarité	$A : B \prec C$	B précède C
Implication	$A : B \Rightarrow C$	si $\exists B$, alors $\exists C$
Exclusion	$A : B \not\wedge C$	pas de B et C simultanément
Constituance	$A : S?$	les descendants $\in S$
Adjacence	$A : B \oplus C$	A est adjacent à C
Dépendance	$A : B \rightsquigarrow C$	B dépend de C

Exemple : le SA

<i>Constituance</i>	AP: {AdP, A, VPinf, PP, Ssub, AP, NP}?
<i>Linéarité</i>	AP: A \prec {VPinf, Ssub, PP, NP, AP} AP: AdP \prec {A, Ssub, PP} AP: AP \prec {A, AdP} AP: PP \prec {Ssub}
<i>Dépendance</i>	AP: {AdP, VPinf, PP, Ssub, NP} \rightsquigarrow A
<i>Unicité</i>	AP: {A, VPinf, Ssub}!
<i>Obligation</i>	AP: Δ A
<i>Exclusion</i>	AP: VPinf $\not\Leftarrow$ {PP, Ssub}

Exemple de graphe de contraintes



Graphe des propriétés satisfaites pour "L'industrie est capable d'investir."

Des contraintes multidimensionnelles

- Relation contour/groupe prosodique:

$$\text{Min_Rising}_{[proso]} : \{ \text{AP}_{[proso]} \}$$

Les syntagmes accentuels se terminent par des contours à montée mineure

- Relation groupe syntaxique/contour :

$$\text{NV}_{[synt]} : \{ \text{V}_{[synt]}, \text{Min_Rising}_{[proso]} \}$$

Les noyaux verbaux sont réalisés par des contours à montée mineure

- Backchannels :

$$\begin{aligned} \text{IP}_{[proso]} &\prec \text{VocBC}_{[disc]} \\ \text{VocBC}_{[disc]} &\prec \{ \text{RF}, \text{RT} \}_{[proso]} \\ \text{VocBC}_{[disc]} &\prec \text{TCU_f}_{[disc]} \end{aligned}$$

Les backchannels apparaissent après des IP, réalisés par des contours montant-descendants ou montées terminales, et après des tours conversationnels complets

Des contraintes multidimensionnelles

- Relation contour/groupe prosodique:

$$\text{Min_Rising}_{[proso]} : \{ \text{AP}_{[proso]} \}$$

Les syntagmes accentuels se terminent par des contours à montée mineure

- Relation groupe syntaxique/contour :

$$\text{NV}_{[synt]} : \{ \text{V}_{[synt]}, \text{Min_Rising}_{[proso]} \}$$

Les noyaux verbaux sont réalisés par des contours à montée mineure

- Backchannels :

$$\begin{aligned} \text{IP}_{[proso]} &\prec \text{VocBC}_{[disc]} \\ \text{VocBC}_{[disc]} &\prec \{ \text{RF}, \text{RT} \}_{[proso]} \\ \text{VocBC}_{[disc]} &\prec \text{TCU_f}_{[disc]} \end{aligned}$$

Les backchannels apparaissent après des IP, réalisés par des contours montant-descendants ou montées terminales, et après des tours conversationnels complets

Des contraintes multidimensionnelles

- Relation contour/groupe prosodique:

$$\text{Min_Rising}_{[proso]} : \{ \text{AP}_{[proso]} \}$$

Les syntagmes accentuels se terminent par des contours à montée mineure

- Relation groupe syntaxique/contour :

$$\text{NV}_{[synt]} : \{ \text{V}_{[synt]}, \text{Min_Rising}_{[proso]} \}$$

Les noyaux verbaux sont réalisés par des contours à montée mineure

- Backchannels :

$$\begin{aligned} \text{IP}_{[proso]} &\prec \text{VocBC}_{[disc]} \\ \text{VocBC}_{[disc]} &\prec \{ \text{RF}, \text{RT} \}_{[proso]} \\ \text{VocBC}_{[disc]} &\prec \text{TCU_f}_{[disc]} \end{aligned}$$

Les backchannels apparaissent après des IP, réalisés par des contours montant-descendants ou montées terminales, et après des tours conversationnels complets

L'analyse en GP

- 1 **Input** : ensemble \mathcal{C} des catégorisations possibles
- 2 Génération de l'ensemble \mathcal{A} des **affectations** (ensemble des sous-ensembles de \mathcal{C})
- 3 Construction des **caractérisations** de \mathcal{A} : évaluation des propriétés décrivant les relations entre catégories pour chaque affectation
- 4 Inférence de **nouvelles catégories** : si une caractérisation correspond à une catégorie, on l'ajoute à \mathcal{C}
- 5 Recherche d'une **solution** : partition de \mathcal{A}
 - qui couvre la totalité de l'input
 - dans lequel chaque mot correspond à une étiquette unique
 - qui satisfait un nombre suffisant de propriétés

L'analyse en GP

- 1 **Input** : ensemble \mathcal{C} des catégorisations possibles
- 2 Génération de l'ensemble \mathcal{A} des **affectations** (ensemble des sous-ensembles de \mathcal{C})
- 3 Construction des **caractérisations** de \mathcal{A} : évaluation des propriétés décrivant les relations entre catégories pour chaque affectation
- 4 Inférence de **nouvelles catégories** : si une caractérisation correspond à une catégorie, on l'ajoute à \mathcal{C}
- 5 Recherche d'une **solution** : partition de \mathcal{A}
 - qui couvre la totalité de l'input
 - dans lequel chaque mot correspond à une étiquette unique
 - qui satisfait un nombre suffisant de propriétés

L'analyse en GP

- 1 **Input** : ensemble \mathcal{C} des catégorisations possibles
- 2 Génération de l'ensemble \mathcal{A} des **affectations** (ensemble des sous-ensembles de \mathcal{C})
- 3 Construction des **caractérisations** de \mathcal{A} : évaluation des propriétés décrivant les relations entre catégories pour chaque affectation
- 4 Inférence de **nouvelles catégories** : si une caractérisation correspond à une catégorie, on l'ajoute à \mathcal{C}
- 5 Recherche d'une **solution** : partition de \mathcal{A}
 - qui couvre la totalité de l'input
 - dans lequel chaque mot correspond à une étiquette unique
 - qui satisfait un nombre suffisant de propriétés

L'analyse en GP

- 1 **Input** : ensemble \mathcal{C} des catégorisations possibles
- 2 Génération de l'ensemble \mathcal{A} des **affectations** (ensemble des sous-ensembles de \mathcal{C})
- 3 Construction des **caractérisations** de \mathcal{A} : évaluation des propriétés décrivant les relations entre catégories pour chaque affectation
- 4 Inférence de **nouvelles catégories** : si une caractérisation correspond à une catégorie, on l'ajoute à \mathcal{C}
- 5 Recherche d'une **solution** : partition de \mathcal{A}
 - qui couvre la totalité de l'input
 - dans lequel chaque mot correspond à une étiquette unique
 - qui satisfait un nombre suffisant de propriétés

L'analyse en GP

- 1 **Input** : ensemble \mathcal{C} des catégorisations possibles
- 2 Génération de l'ensemble \mathcal{A} des **affectations** (ensemble des sous-ensembles de \mathcal{C})
- 3 Construction des **caractérisations** de \mathcal{A} : évaluation des propriétés décrivant les relations entre catégories pour chaque affectation
- 4 Inférence de **nouvelles catégories** : si une caractérisation correspond à une catégorie, on l'ajoute à \mathcal{C}
- 5 Recherche d'une **solution** : partition de \mathcal{A}
 - qui couvre la totalité de l'input
 - dans lequel chaque mot correspond à une étiquette unique
 - qui satisfait un nombre suffisant de propriétés

Exemple

<i>La</i>	<i>principale</i>	<i>activité</i>	<i>du</i>	<i>domaine</i>	<i>...</i>
Det	Adj	N	Prep	N	
Pro	N				
N					

- Affectations $\mathcal{A} = \{\text{Det}\}; \{\text{Adj}\}; \{\text{Det}, \text{Adj}\}; \{\text{Det}, \text{Adj}, \text{N}\}; \{\text{Pro}, \text{Adj}\}, \{\text{N}, \text{Prep}, \text{N}\}, \text{etc.}$
- Caractérisation de $\{\text{Det}\} : \mathcal{C}^+ = \{NP : \text{Det}!\}$ $\mathcal{C}^- = \{NP : \text{Det} \Rightarrow N\}$
Caractérisation de $\{\text{Adj}\} : \mathcal{C}^+ = \{AP : \text{Adj}!; AP : \Delta\text{Adj}\}$
- Inférence de AP

Facteurs de complexité

- Complexité classique des formalismes DI/PL : dépend du nombre de contraintes linéaires
- Taille de l'ensemble des catégories :
 - Croissant au fur et à mesure de l'analyse
 - Ensemble potentiellement grand : prise en compte de toutes les dimensions
- Taille de l'espace de recherche : ensemble des affectations générées
- Recherche d'une partition couvrant l'input
- Evaluation des propriétés : fonction de la taille de la grammaire

Données initiales : treebank enrichi

```
<category label="NP" features="NP:SUJ" node_index="0:16:1">
  <category label="Det" features="Da-ms----" node_index="0:16:2" form="Le"
  <category label="Noun" features="Ncms--" node_index="0:16:3" form="text"
  <category label="AP" features="AP" node_index="0:16:4">
    <category label="Adj" features="Af-ms-" node_index="0:16:5" form="mé"
  </category>
  <characterization>
    <property type="lin" source="0:16:2" target="0:16:3" sat="p"/>
    <property type="lin" source="0:16:2" target="0:16:4" sat="p"/>
    <property type="req" source="0:16:3" target="0:16:2" sat="p"/>
    <property type="dep" source="0:16:2" target="0:16:3" sat="p"/>
    <property type="dep" source="0:16:4" target="0:16:3" sat="p"/>
    <property type="oblig" source="0:16:1" target="0:16:3" sat="p"/>
    <property type="uniq" source="0:16:1" target="0:16:2" sat="p"/>
  </characterization>
</category>
```

Affectations caractérisées

Propriétés	Réalizations				
	{Det, N}	{Det, N, AP}	{Det, N, PP}	{Det, N, AP, PP}	{N, NP}
[1] Det \prec N	x	x	x	x	
[2] Det \prec AP		x		x	
[3] Det \prec PP			x	x	
[4] Det \prec NP					x
[5] N \prec PP			x	x	
[6] Oblig(N)	x	x	x	x	x
[7] Det \rightsquigarrow N	x	x	x	x	
[8] PP \rightsquigarrow N			x	x	
[9] AP \rightsquigarrow N		x		x	
[10] Det \Rightarrow N	x	x	x	x	x
[11] PP \Rightarrow N			x	x	
[12] AP \Rightarrow N		x		x	
[13] Unic(Det)	x	x	x	x	
[14] Unic(N)	x	x	x	x	
[15] N $\not\Rightarrow$ Pro	x	x	x	x	x
[16] N $\not\Rightarrow$ Npro	x	x	x	x	x
[17] AP $\not\Rightarrow$ Pro		x		x	
...					

Prédiction

- En CFG : étant donné un ensemble de catégories (contexte), on prédit une transition (nouvelle catégorie)
- En GP : étant donné un ensemble de contraintes, on peut :
 - prédire une transition
 - prédire une caractérisation
 - prédire une projection (inférer une catégorie supérieure)

Contrôle dans le choix de l'affectation

Affectation	{Det, N, AP}	{Det, N, PP}
Distribution	0.43	0.2
Caractérisation	Det \prec N, Det \prec AP, Δ N, Det \rightsquigarrow N, AP \rightsquigarrow N, Det \Rightarrow N, AP \Rightarrow N, Det!, N!, N $\not\Leftarrow$ Pro, N $\not\Leftarrow$ Npro, AP $\not\Leftarrow$ Pro	Det \prec N, Det \prec PP, N \prec PP, Δ N, Det \rightsquigarrow N, PP \rightsquigarrow N, Det \Rightarrow N, PP \Rightarrow N, Det!, N!, N $\not\Leftarrow$ Pro, N $\not\Leftarrow$ Npro, PP $\not\Leftarrow$ Ssub

Caractérisations multidimensionnelles

Exemple :

Marie	je	la	supporte	pas
N	Clit	Pro	V	Adv
NP	S			

Dislocation	
Disl:	NP?, S?, L*H?
Disl:	NP \prec S
Disl:	NP \oplus S
Disl:	NP \rightsquigarrow S//Clit
Disl:	NP \approx L*H

Vocatif	
Voc:	NP?, S?, H*L?
Voc:	NP \prec S
Voc:	NP \oplus S
Voc:	NP \approx H*L

Prédiction de projection

Distribution :

- $\{NP?, S?, NP \prec S, NP \approx L^*H\} = 0.7$
- $\{NP?, S?, NP \prec S, NP \approx H^*L\} = 0.3$

En rencontrant $\{NP?, S?, NP \prec S\}$, on prédit comme structure la plus probable la projection de la dislocation en même temps que sa caractérisation complète:

Disl: $\{NP?, S?, L^*H?, NP \oplus S, NP \prec S, NP \rightsquigarrow S//Clit, NP \approx L^*H, \dots\}$

Processus

$$\{m_1 \dots m_n\}$$



$$\begin{bmatrix} c_{[1,1]} \\ \dots \\ c_{[1,u]} \end{bmatrix} \dots \begin{bmatrix} c_{[n,1]} \\ \dots \\ c_{[n,b]} \end{bmatrix}$$



$$\{c_{[1,t]} \dots c_{[i,u]}\} \dots \{c_{[j,v]} \dots c_{[n,w]}\}$$



$$\begin{bmatrix} \text{LABEL } l_1 \\ \text{AFFECT } \{c_{[1,t]} \dots c_{[i,u]}\} \\ \text{CARAC } \{p_1, \dots, p_l\} \end{bmatrix} \dots \begin{bmatrix} \text{LABEL } l_v \\ \text{AFFECT } \{c_{[j,t]} \dots c_{[n,u]}\} \\ \text{CARAC } \{p_1, \dots, p_m\} \end{bmatrix}$$

Catégories probabilisées

Affectations probabilisées

Caractérisations probabilisées

Conclusion

- Données observables : treebanks hybrides
- Contrôle
 - Choix des affectations
 - Calcul des caractérisations
- Intérêts
 - Informations partielles
 - Informations sous-spécifiées
 - Données non canoniques