

# Discovery of association rules between syntactic variables

Seminar in Methodology and Statistics, Groningen, 23 May 2007, Marco René Spruit  
<http://www.mre.nl/naar/nl/mre/stevenker/spruit/naar/collocat.htm>

## The big picture

- Generative syntax and functional typology share a primary interest in understanding the structural similarities and differences between language varieties
  - Ultimate goal: to characterise the superficial structural diversity of all language varieties as particular settings of relatively few parametric patterns
- This contribution: A computational method to automatically discover syntactic variable associations

4/22

## SAND1 domains

1. Complementisers
  - 't iijkt wel **of** er iemand in de tuin staat.  
*"It looks AFFIRM if there someone in the garden stands"*
2. Subject pronouns
  - Ze geloof dat **jif** eerder thuis bent dan ik.  
*"She believes that YOU earlier home are than I"*
3. Subject doubling
  - As **ge** **gij** gezond leeft, leef **de** **gij** langer.  
*"If YOU weak YOU strong healthily live, live YOU weak YOU strong longer"*
4. Reflexive and reciprocal pronouns
  - Jan herinnert **zich** dat verhaal wel.  
*"John remembers himself that story AFFIRM"*
5. Fronting
  - Dat is de man **die** het verhaal heeft verteld.  
*"that is the man who the story has told"*

9/22

## Research context

- The Determinants of Dialectal Variation project (DDV)
  - <http://dialectometry.net>
  - University of Groningen: information science
    - John Nerbonne
    - Wilbert Heeringa
  - Meertens Instituut: syntactic theory
    - Hans Bennis
    - Sjef Barbiers
  - "What are the determinants of dialectal variation?"

2/22

## Syntactic variation data

- Syntactic Atlas of the Dutch Dialects (SAND)
  - 267 Dutch dialects
  - SAND1: [Barbiers et al. 2005]
    - Complementisers, Subject pronouns, Subject doubling, Reflexive and reciprocal pronouns, Fronting
    - 106 syntactic contexts, 485 variables
  - SAND2: [Barbiers et al. 2007]
    - Verbal clusters, Cluster interruption, Morphosyntactic variation, Negative particle, Negative concord and quantification
    - 65 syntactic contexts, 274 variables (*incomplete*)

5/22

## Data mining the SAND

- Knowledge Discovery in Databases (KDD)
  - "the science of extracting useful information from large data sets or databases" (Hand et al., 2001)
  - An umbrella term for techniques like association rules, decision trees, neural networks, ...
  - Association rule mining: A → C
    - A: predicting attribute value(s) ("antecedent")
    - C: predicted class ("consequent")
  - Based on proportional overlap
    - Geographical co-occurrences of variables

11/22

## Syntactic variation & dialectometry

- Language variation dimensions
  - { Macro, **Micro** }
  - { Pronunciation, Lexis, Morphology, **Syntax** }
  - { External, **Internal** }
  - { Time, **Space** }
  - { Qualitative, **Quantitative** }
- Research questions
  - i. How can relevant associations between syntactic variables be discovered?
  - ii. What are interesting associations between syntactic variables?

3/22

## Dutch language area

- Distribution of the 267 Dutch dialects in the SAND
- The provinces in the Dutch language area

6/22

## Sample variables

- A. "Complementiser of comparative if-clause" (14b)
 

't iijkt wel	<b>of</b>	iemand	in de tuin	staat.
it looks [affirm]	if	that there someone	in the garden	stands
- B. "Subject doubling 2 singular" (54a)
 

Ge	geloof	gij	zeker	niet dat hij sterker is as	<b>ge</b>	<b>gij</b> .
you weak	believe	you strong	certainly not	that he stronger is than you weak	you strong	
- C. "Weak reflexive pronoun as object of inherent reflexive verb" (68a)
 

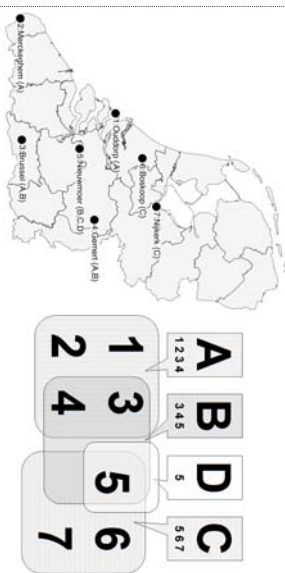
Jan	herinnert	<b>zich</b>	dat	verhaal	wel.
John	remembers	his own	that	story	[affirmative]
- D. "Short subject relative, complementiser following relative pronoun" (84a)
 

Dat is de man	<b>die</b>	<b>dat</b>	het verhaal	verteld heeft.
that is the man	who	that	the story	told has

12/22

## Sample data illustration

- Example: 4 variables (A-D) in 7 locations (1-7)



## Evaluation factors of rule quality

- **Accuracy:**  $IA \& CI / IA I$   
How often is the rule correct?  
- varA → varB:  $(A \cap B / A) * 100 = 2/4 * 100 = 50\%$
- **Coverage:**  $IA I$   
How often does the rule apply?  
- varA → varB:  $A / N * 100 = 4/7 * 100 = 57\%$
- **Completeness:**  $IA \& CI / ICI$   
How much of the target class does the rule cover?  
- varA → varB:  $(A \cap B / B) * 100 = 2/3 * 100 = 66\%$
- **Interestingness:**  $IA \& CI - IA ICI / N$   
Integrates the three factors above into one value...  
- varA → varB:  $(A \cap B) - (A * B / N) = 2 - (4 * 3 / 7) = 0.28$

## Sample data results

- The 8 highest ranked association rules:

#	Antecedent → Consequent	Interestingness	Complexity	Accuracy	Coverage	Completeness
1.	B → A ∨ D	0.86	1	100	42	60
2.	A ∨ D → B	0.86	1	60	71	100
3.	D → B	0.57	0	100	14	33
4.	D → C	0.57	0	100	14	33
5.	B → D	0.57	0	33	42	100
6.	C → D	0.57	0	33	42	100
7.	B → A	0.29	0	66	42	50
8.	A → B	0.29	0	50	57	66

## No. 1 association rule in SAND1

A1n1: p46a:g-lieden (Subject pronouns 2 plural, strong forms)  
We geloven dat **g-lieden** niet zo slim zijn als wij.  
We believe that *you plural/strong* not so smart are as we.  
We believe that you are not as smart as we are.'

Cons: p38b:gij/gie (Subject pronouns 2 singular, strong forms)  
Ze geloof dat **gij/gie** eerder thuis bent dan ik.  
she believes that *you singular/strong* earlier home are than I.  
She thinks that you'll be home sooner than me.'

Stat: Rank=1, Combination=10,321, Interestings=58,38, Accuracy=99%, Coverage=39%, Completeness=89%, Complexity=0, A-Locations=105, C-Locations=116, AC-Overlap=104, AC-Disjunction=117

Interp: The plural pronoun 'g-lieden' belongs to the same paradigm as the singular pronoun 'gij'.

## More associated rules

- We geloven dat g-lieden niet zo slim zijn als wij.  
*'we believe that you plural/strong not so smart are as we'*
- a) Ze geloof dat gij/gie eerder thuis bent dan ik.  
*'she believes that you earlier home are than I'*
- b) Ik denk da Marie hem zal moeten roepen.  
*'I think that Mary him will must call'*
- c) U [niet-beleefdh] geloof dat Lisa even mooi is als Anna.  
*you [non-polite/cj] believe that Lisa as beautiful is as Anna'*
- d) Fons zag een slang naast hem.  
*'Fons saw a snake next to him'*
- e) Erik liet mij voor hem werken.  
*'Erik let me for him work'*
- f) De jongen wie/die zijn moeder yesterday remarried is.  
*'the boy who/that his mother yesterday remarried is'*

## A higher complexity rule

- "If either antecedent variable A1 or A2 occurs in a dialect, then syntactic variable C also occurs"

A1: p46b:julle(n)/julle (Subject pronouns 2 plural, strong forms, complex)  
We geloven dat **julle(n)/julle** niet zo slim zijn als wij.  
we believe that *you plural/strong* not so smart are as we.  
'We believe that you are not as smart as we are.'

A2: p46b:julder/jelder (Subject pronouns 2 plural, strong forms, complex)  
We geloven dat **julder/jelder** niet zo slim zijn als wij.

C: p46a:j-lieden-compositum (Subject pronouns 2 plural, strong forms)  
We geloven dat **lieden** niet zo slim zijn als wij.

Int: The infrequent pronoun 'julder/jelder' perfects the implicational association of the frequent 'julle(n)/julle' variant with the pronoun 'j-lieden'.

## Some conclusions

1. Association rule mining technique based on proportional overlap: *It works.*
  - Facilitates identification, validation and exploration of variable relationships
2. Reveals the existence of many potentially interesting associations within SAND1
3. Shows considerable overlaps between the geographical distributions of syntactic variable pairs
4. Results strongly indicate that many more potentially interesting associations between syntactic variables are likely to be uncovered

## Implicational chain of rules

1/4: d54a:after v (Subject doubling 2 singular)  
As gij gezond leeft, leef-**de** **gij** langer.  
If *you sing* healthily live, live- *you plural/weak* *you sing/strong* longer

2/4: d55a:after v (Subject doubling 2 plural)  
As gilder gezond leeft, leef-**de** **gilder** langer.  
If *you plural* healthily live, live- *you plural/weak* *you plural/strong* longer

3/4: p46a:g-lieden (Subject pronouns 2 plural, strong forms)  
We geloven dat **g-lieden** niet zo slim zijn als wij.  
we believe that *you plural/strong* not so smart are as we.

4/4: p38b:gij/gie (Subject pronouns 2 singular, strong forms)  
Ze geloof dat **gij/gie** eerder thuis bent dan ik.  
she believes that *you singular/strong* earlier home are than I

## Discussion & future research

- Incorporate exception rules
- Alternative measures of interestingness / incorporation of additional rule quality evaluation factors (Surprisingness, ...)
- Adding more data (SAND2)
  - Phonological data: discover potential associations between variables among *linguistic levels*
- Refine dialect area detection
- Comparison with methods such as Cramér's V and correspondence analysis