

# Quantity and quality in linguistic variation

## The case of verb clusters

Jeroen van Craenenbroeck

KU Leuven/CRISSP

Séminaire

INALCO – SeDyL

Paris, February 6, 2015

This talk in one  
slide

Introduction

The data: dialect  
Dutch verb clusters

Research questions

Theoretical  
background:  
dialectometry

Methodology:  
reverse  
dialectometry

Results

Main conclusion

Extensions

References

# This talk in one slide

- ▶ **main topic:** interaction between formal-theoretical and quantitative-statistical linguistics
- ▶ **starting point:** the massive amount of variation attested in Dutch verb clusters necessitates a collaboration between formal and quantitative approaches
- ▶ **traditional dialectometry** measures (dis)similarities between dialect locations based on their linguistic profile
- ▶ **reverse dialectometry** measures (dis)similarities between *linguistic constructions* based on their geographical distribution and maps these results against formal-theoretical parameters
- ▶ **result:** a method that can detect and identify grammatical parameters in a large and highly varied data set

Quantity and  
quality in linguistic  
variation

Jeroen van  
Craenenbroeck

This talk in one  
slide

Introduction

The data: dialect  
Dutch verb clusters

Research questions

Theoretical  
background:  
dialectometry

Methodology:  
reverse  
dialectometry

Results

Main conclusion

Extensions

References

# Introduction

- ▶ quantitative and qualitative linguistics typically represent two non-communicating, non-intersecting subfields
- ▶ **qualitative linguistics:**
  - ▶ theory-driven
  - ▶ works with small and/or simplified data sets
  - ▶ makes little or no use of mathematical or statistical methods
- ▶ **quantitative linguistics:**
  - ▶ data-driven
  - ▶ works with (very) large and highly varied data sets
  - ▶ makes extensive use of mathematical or statistical methods

Quantity and  
quality in linguistic  
variation

Jeroen van  
Craenenbroeck

This talk in one  
slide

Introduction

The data: dialect  
Dutch verb clusters

Research questions

Theoretical  
background:  
dialectometry

Methodology:  
reverse  
dialectometry

Results

Main conclusion

Extensions

References

- ▶ **this talk** explores the possibility of a collaboration between the two approaches:
  - ▶ to what extent can a quantitative analysis of large datasets lead to new theoretical insights?
  - ▶ to what extent can theoretical analyses guide and inform quantitative analyses of language data?
- ▶ **proposal:**
  - ▶ quantitative analyses can be used to test, evaluate, and compare hypotheses from the theoretical literature
  - ▶ theoretical analyses can be used to interpret and narrow down the space of hypotheses generated by quantitative analyses
- ▶ **main topic for today:** word order variation in dialect Dutch verb clusters:
  - ▶ lots of data available and lots of interdialectal variation
  - ▶ extensively researched in the formal-theoretical (esp. generative) literature for at least four decades

## The data: dialect Dutch verb clusters

- ▶ in Dutch (like in many Germanic languages) verbs tend to group together at the right edge of the (embedded) clause:

(1) dat hij gisteren tijdens de les **gelachen** heeft.  
 that he yesterday during the class laughed has  
 'that he laughed yesterday during class.' (21)

The data: dialect  
Dutch verb clusters

- ▶ moreover, such verbal clusters typically show a certain degree of freedom in their word order:

(2) dat hij gisteren tijdens de les heeft gelachen.  
that he yesterday during the class had laughed  
'that he laughed yesterday during class.' (12)

## Extensions

## References

- ▶ this word order freedom is typically a source of interdialectal variation:

### (3) Ferwerd Dutch

- a. dasto it ook net **zien** **meist**.  
that.you it also not see may  
'that you're also not allowed to see it.' (✓**21**)
- b. \*dasto it ook net **meist** **zien**.  
that.you it also not may see  
'that you're also not allowed to see it.' (\***12**)



- ▶ this word order freedom is typically a source of interdialectal variation:

#### (4) Gendringen Dutch

- a. dat ee et ook nie **zien** **mag**.  
that you it also not see may  
'that you're also not allowed to see it.' (✓**21**)
- b. dat ee et ook nie **mag** **zien**.  
that you it also not may see  
'that you're also not allowed to see it.' (✓**12**)

● Gendringen (1)



- this word order freedom is typically a source of interdialectal variation:

## (5) Poelkapelle Dutch

- a. \*dajtgie ook nie **zien** **meug**.  
that.it.you also not see may  
'that you're also not allowed to see it.' (\*21)
- b. dajtgie ook nie **meug** **zien**.  
that.it.you also not may see  
'that you're also not allowed to see it.' (✓12)





- ▶ and the more complex the verbal cluster, the more variation there is: in verbal clusters consisting of two modal auxiliaries and one main verb, out of the six orders that are theoretically possible, four are attested in Dutch dialects:

- (6) Ik vind dat iedereen moet kunnen zwemmen.  
I find that everyone must can swim  
'I think everyone should be able to swim.' (✓123)
- (7) a. ...dat iedereen moet zwemmen kunnen. (✓132)  
b. ...dat iedereen zwemmen moet kunnen. (✓312)  
c. ...dat iedereen zwemmen kunnen moet. (✓321)  
d. \*...dat iedereen kunnen zwemmen moet. (\*231)  
e. \*...dat iedereen kunnen moet zwemmen. (\*213)

- ▶ but once again, it is not the case that each of the four allowed orders is attested in all dialects:

(8) *Midsland Dutch*

- a. \*dat elkeen mot kanne zwemme.  
that everyone must can swim  
'that everyone should be able to swim.' (\*123)
- b. dat elkeen mot zwemme kanne. (✓132)
- c. \*dat elkeen zwemme mot kanne. (\*312)
- d. dat elkeen zwemme kanne mot. (✓321)
- e. \*dat elkeen kanne zwemme mot. (\*231)
- f. \*dat elkeen kanne mot zwemme. (\*213)

- but once again, it is not the case that each of the four allowed orders is attested in all dialects:

(9) *Langelo Dutch*

- a. dat iedereen mot kunnen zwemmen.  
that everyone must can swim  
'that everyone should be able to swim.' (✓123)
- b. \*dat iedereen mot zwemmen kunnen. (\*132)
- c. dat iedereen zwemmen mot kunnen. (✓312)
- d. \*dat iedereen zwemmen kunnen mot. (\*321)
- e. \*dat iedereen kunnen zwemmen mot. (\*231)
- f. \*dat iedereen kunnen mot zwemmen. (\*213)

- more generally, the four possible cluster orders yield a total of 16 possible combinations, of which 12 are attested in Dutch dialects:

<b>example dialect</b>	<b>123</b>	<b>132</b>	<b>321</b>	<b>312</b>
Beetgum	✓	✓	✓	✓
Hippolytushoef	✓	✓	✓	*
Warffum	✓	✓	*	*
Oosterend	✓	*	*	*
Schermerhorn	✓	✓	*	✓
Visvliet	✓	*	✓	✓
Kollum	✓	*	✓	*
Langelo	✓	*	*	✓
Midsland	*	✓	✓	*
Lies	*	*	✓	*
Bakkeveen	*	*	✓	✓
Waskemeer	*	✓	*	*

This talk in one  
slide

Introduction

The data: dialect  
Dutch verb clusters

Research questions

Theoretical  
background:  
dialectometryMethodology:  
reverse  
dialectometry

Results

Main conclusion

Extensions

References

- ▶ in order to get a more complete picture of the variation, we can look at the results from the SAND-project:
  - ▶ Syntactic Atlas of the Dutch Dialects (2000–2004)
  - ▶ dialect interviews in 267 dialect locations in Belgium, France, and the Netherlands
- ▶ the SAND-questionnaire contained eight questions on word order in verb clusters for a total of 31 cluster orders
- ▶ if we map, for each of the 267 SAND-dialects, which dialect has which combination of cluster orders, we find 137 different combinations of verb cluster orders
- ▶ in other words, there are 137 different types of dialects when it comes to word order in verbal clusters

# Research questions

- ▶ this state of affairs raises (at least) two fundamental questions for grammatical theory:
  1. to what extent is this variation due to grammar-internal properties, and what part of it is grammar-external?
  2. how can we draw the line between the two?
- ▶ one extreme position: Barbiers (2005): the grammar rules out 231 and 213 in MOD-MOD-V-cluster, but all other orders are freely available to all speakers; the choice between them is determined by sociolinguistic factors (geographical and social norms, register, context, ...)
- ▶ **in this talk** I use quantitative-statistical methods to identify three grammatical (micro)parameters that together are responsible for the bulk of the variation

Quantity and  
quality in linguistic  
variation

Jeroen van  
Craenenbroeck

This talk in one  
slide

Introduction

The data: dialect  
Dutch verb clusters

Research questions

Theoretical  
background:  
dialectometry

Methodology:  
reverse  
dialectometry

Results

Main conclusion

Extensions

References

# Theoretical background: dialectometry

- ▶ **dialectometry** is a subdiscipline of linguistics that uses computational and quantitative techniques in dialectology (Nerbonne and Kretzschmar Jr., 2013)
- ▶ this section presents a prototypical dialectometric analysis, which will serve as a stepping stone for the actual analysis in the next section
- ▶ **starting point:** the raw data from 13 SAND-maps:
  - ▶ 4 about two-verb clusters ( $3 \times$  AUXILIARY-PARTICIPLE,  $1 \times$  MODAL-INFINITIVE)
  - ▶ 4 about three-verb clusters (MODAL-MODAL-INFINITIVE, MODAL-AUXILIARY-PARTICIPLE, AUXILIARY-AUXILIARY-INFINITIVE, AUXILIARY-MODAL-INFINITIVE)
  - ▶ 3 about particle placement inside the cluster
  - ▶ 2 about morphology of the past participle
- ▶ for a total of 67 linguistic variables in 267 locations

- ▶ this yields a  $267 \times 67$  matrix with one row per location and one column per linguistic variable, i.e. locations = individuals and linguistic phenomena = variables



	AUX1(be.sg)-PART2	PART2-AUX1(be.sg)	AUX1(have.sg)-PART2	PART2-AUX1(have.sg)	AUX1(have
Midslân / Midslân	no	yes	no	yes	
Lies	no	yes	no	yes	
West-Terschelling	no	yes	no	yes	
Oosterend	NA	NA	no	yes	
Hollum	no	yes	NA	NA	
Schiermonnikoog	no	yes	no	yes	
Ferwerd / Ferwert	no	yes	no	yes	
Anjum / Eanjum	no	yes	no	yes	
Kollum	no	yes	no	yes	
Visvliet	no	yes	no	yes	
Oosterbierum / Ea	no	yes	no	yes	
Beetgum / Bitgum	no	yes	NA	NA	
Bergum / Burgum	no	yes	no	yes	
Jorwerd / Jorwert	no	yes	NA	NA	
Bakkeveen / Bakke	no	yes	no	yes	
Waskemeer / De V	no	yes	no	yes	
Kloosterburen	no	yes	no	yes	
Warffum	no	yes	no	yes	
Leermens	no	yes	no	yes	
Groningen	no	yes	yes	no	
Nieuw-Scheemda	NA	NA	no	yes	
Langelo	no	yes	no	yes	

- ▶ this yields a  $267 \times 67$  matrix with one row per location and one column per linguistic variable, i.e. locations = individuals and linguistic phenomena = variables
- ▶ step 1: convert the table into a  $267 \times 267$  (symmetric) distance matrix, whereby for each pair of locations a distance between them is calculated based on the linguistic features they share

	Midslan	Lies	West-Ter	Oosteren	Hollum	Schiermc	Ferwerd	Anjum /	Kollum	Visvliet	Oosterbie	Beet
Midslan / M	0,000	0,500	0,333	0,706	0,250	0,647	0,357	0,250	0,611	0,650	0,533	0,
Lies	0,500	0,000	0,444	0,750	0,588	0,375	0,471	0,563	0,444	0,444	0,632	0,
West-Tersch	0,333	0,444	0,000	0,789	0,429	0,667	0,286	0,429	0,632	0,600	0,500	0,
Oosterend	0,706	0,750	0,789	0,000	0,706	0,765	0,737	0,538	0,563	0,600	0,600	0,
Hollum	0,250	0,588	0,429	0,706	0,000	0,667	0,167	0,000	0,625	0,714	0,462	0,
Schiermonnik	0,647	0,375	0,667	0,765	0,667	0,000	0,625	0,667	0,400	0,556	0,706	0,
Ferwerd / Fer	0,357	0,471	0,286	0,737	0,167	0,625	0,000	0,182	0,588	0,682	0,308	0,
Anjum / Eanj	0,250	0,563	0,429	0,538	0,000	0,667	0,182	0,000	0,571	0,625	0,417	0,
Kollum	0,611	0,444	0,632	0,563	0,625	0,400	0,588	0,571	0,000	0,353	0,625	0,
Visvliet	0,650	0,444	0,600	0,600	0,714	0,556	0,682	0,625	0,353	0,000	0,588	0,
Oosterbierum	0,533	0,632	0,500	0,600	0,462	0,706	0,308	0,417	0,625	0,588	0,000	0,
Beetgum / Bit	0,545	0,714	0,500	0,727	0,500	0,750	0,333	0,556	0,643	0,500	0,167	0,
Bergum / Bur	0,500	0,500	0,429	0,813	0,500	0,571	0,333	0,500	0,429	0,667	0,571	0,
Jorwerd / Jor	0,692	0,667	0,583	0,846	0,545	0,667	0,400	0,600	0,571	0,692	0,500	0,
Bakkeveen / t	0,400	0,500	0,438	0,706	0,385	0,563	0,357	0,385	0,438	0,579	0,533	0,
Waskemeer /	0,438	0,526	0,556	0,818	0,500	0,588	0,471	0,533	0,471	0,652	0,588	0,
Kloosterburen	0,500	0,412	0,611	0,810	0,563	0,357	0,529	0,600	0,333	0,636	0,706	0,
Warffum	0,563	0,438	0,667	0,737	0,625	0,429	0,588	0,643	0,400	0,652	0,600	0,
Leermens	0,667	0,652	0,739	0,550	0,773	0,650	0,739	0,722	0,389	0,455	0,667	0,
Groningen	0,714	0,682	0,714	0,636	0,783	0,762	0,800	0,778	0,471	0,476	0,684	0,
Nieuw-Scheer	0,650	0,682	0,650	0,652	0,773	0,762	0,739	0,722	0,556	0,368	0,647	0,
Langelo	0,727	0,524	0,739	0,652	0,792	0,650	0,760	0,647	0,550	0,500	0,700	0,

This talk in one  
slide

Introduction

The data: dialect  
Dutch verb clusters

Research questions

Theoretical  
background:  
dialectometryMethodology:  
reverse  
dialectometry

Results

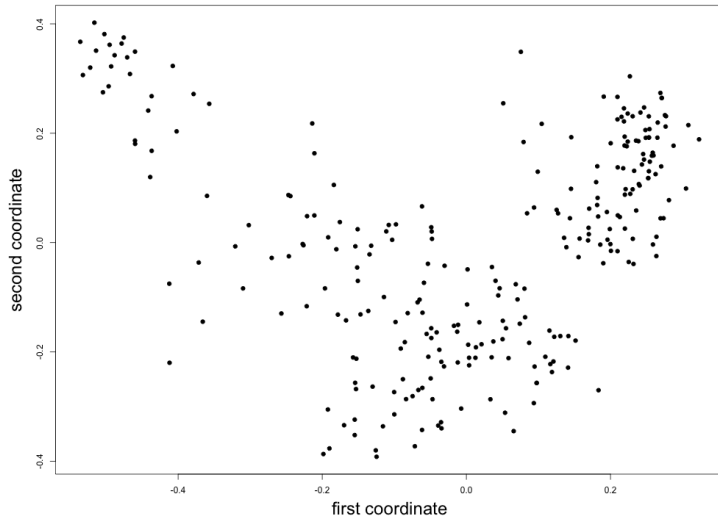
Main conclusion

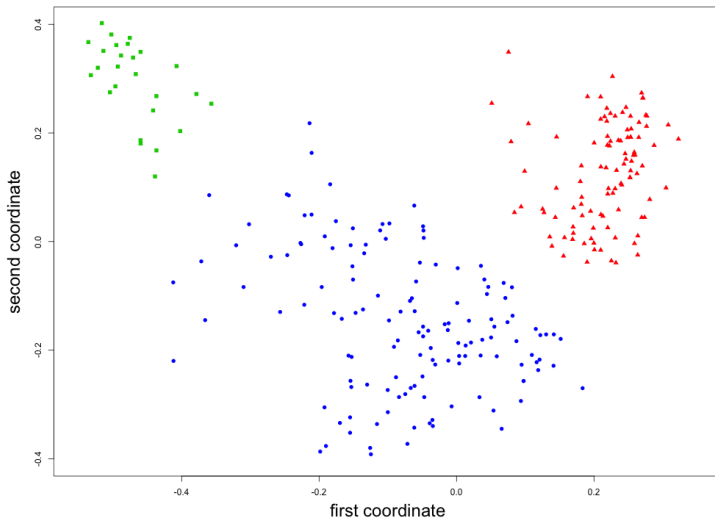
Extensions

References

- ▶ this yields a  $267 \times 67$  matrix with one row per location and one column per linguistic variable, i.e. locations = individuals and linguistic phenomena = variables
- ▶ step 1: convert the table into a  $267 \times 267$  (symmetric) distance matrix, whereby for each pair of locations a distance between them is calculated based on the linguistic features they share
- ▶ step 2: apply multidimensional scaling (MDS) to the distance matrix
- ▶ MDS is a mathematical technique for reducing a multidimensional distance matrix to a low dimensional space in which each point represents an object from the distance matrix, and distances between points represents, as well as possible, dissimilarities between objects (Borg and Groenen, 2005)

## 2-dimensional MDS-representation 67 verb cluster phenomena





This talk in one  
slide

Introduction

The data: dialect  
Dutch verb clusters

Research questions

Theoretical  
background:  
dialectometry

Methodology:  
reverse  
dialectometry

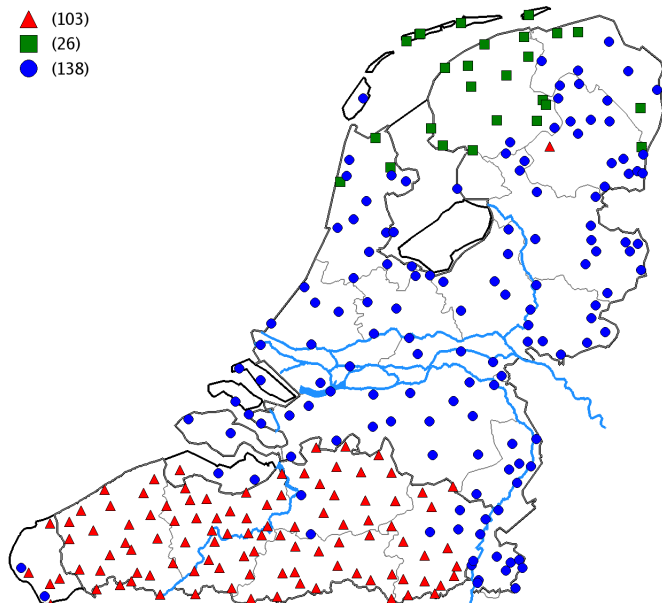
Results

Main conclusion

Extensions

References

- ▶ this yields a  $267 \times 67$  matrix with one row per location and one column per linguistic variable, i.e. locations = individuals and linguistic phenomena = variables
- ▶ step 1: convert the table into a  $267 \times 267$  (symmetric) distance matrix, whereby for each pair of locations a distance between them is calculated based on the linguistic features they share
- ▶ step 2: apply multidimensional scaling (MDS) to the distance matrix
- ▶ step 3: project the data back onto a geographical map





- ▶ **note:** the linguistic variables (i.e. cluster orders) are used to determine the degree of similarity/difference between dialect locations
- ▶ these similarities and differences are then projected back onto a geographical map, which makes it possible to discern dialect regions based on what verb cluster phenomena they possess
- ▶ shortcomings of this approach for our current purposes:
  1. the linguistic constructions themselves play only an indirect role in the outcome of the analysis: we can see when two dialects differ, but we don't see which cluster orders are responsible for this difference and to what extent
  2. there is no link between the data that feed into the quantitative analysis and the formal theoretical literature on verb clusters

# Methodology: reverse dialectometry

- ▶ **proposal:** two changes to the classical dialectometric setup:
  1. cluster orders are *individuals* rather than variables, i.e. instead of calculating differences between dialect locations, we measure differences between linguistic constructions
  2. theoretical analyses of verb cluster orders are decomposed in their constitutive parts, which makes it possible to include them as supplementary variables in the analysis

Quantity and  
quality in linguistic  
variation

Jeroen van  
Craenenbroeck

This talk in one  
slide

Introduction

The data: dialect  
Dutch verb clusters

Research questions

Theoretical  
background:  
dialectometry

Methodology:  
reverse  
dialectometry

Results

Main conclusion

Extensions

References

- ▶ starting point: a  $31 \times 267$  data table whereby each cluster order represents a row and each dialect location a column

	Midsland	Lies	West.Tersch	Oosterend	Hollum	Schiermonni	Ferwerd	Anjum	Kollum
AUX1(be.sg)-PART2	no	no	no	NA	no	no	no	no	no
PART2-AUX1(be.sg)	yes	yes	yes	NA	yes	yes	yes	yes	yes
AUX1(have.sg)-PART2	no	no	no	no	NA	no	no	no	no
PART2-AUX1(have.sg)	yes	yes	yes	yes	NA	yes	yes	yes	yes
AUX1(have.pl)-PART2	no	no	no	no	no	no	no	no	no
PART2-AUX1(have.pl)	yes	yes	yes	yes	yes	yes	yes	yes	yes
MOD1(sg)-INF2	no	no	yes	no	no	no	no	no	no
INF2-MOD1(sg)	yes	yes	yes	yes	yes	yes	yes	yes	yes
MOD2-INF3-MOD1(sg)	no	no	no	no	no	no	no	no	no
MOD1(sg)-MOD2-INF3	no	no	no	yes	no	no	no	no	yes
MOD1(sg)-INF3-MOD2	yes	no	no	no	no	no	no	no	no
INF3-MOD2-MOD1(sg)	yes	yes	yes	no	yes	yes	yes	yes	yes
INF3-MOD1(sg)-MOD2	no	no	no	no	no	no	no	no	no
MOD1(sg)-AUX2(have)-PART3	no	no	no	no	no	no	no	NA	no
MOD1(sg)-PART3-AUX2(have)	no	no	no	no	no	no	no	NA	yes
PART3-MOD1(sg)-AUX2(have)	no	yes	no	yes	no	no	no	NA	yes
PART3-AUX2(have)-MOD1(sg)	yes	yes	yes	no	yes	yes	yes	NA	yes
AUX1(be.sg)-AUX2(go)-INF3	no	no	no	yes	no	no	no	no	NA
AUX1(be.sg)-INF3-AUX2(go)	no	no	no	no	no	no	no	no	NA
AUX2(go)-AUX1(be.sg)-INF3	no	no	no	no	no	yes	no	no	NA
AUX2(go)-INF3-AUX1(be.sg)	no	no	no	no	no	no	no	no	NA
INF3-AUX1(be.sg)-AUX2(go)	no	no	no	no	no	no	no	no	NA
INF3-AUX2(go)-AUX1(be.sg)	yes	yes	yes	no	yes	no	yes	yes	NA
AUX1(have.sg)-MOD2(inf)-INF3	no	no	no	yes	no	no	no	no	no
AUX1(have.sg)-INF3-MOD2(part)	no	no	no	no	no	no	no	no	no
AUX1(have.sg)-INF3-MOD2(inf)	no	no	no	no	no	no	no	no	no
MOD2(inf)-INF3-AUX1(have.sg)	no	no	no	no	no	no	no	no	no
INF3-AUX1(have.sg)-MOD2(inf)	no	no	yes	no	no	no	no	no	no
INF3-AUX1(have.sg)-MOD2(part)	no	no	no	no	no	no	no	no	no
INF3-MOD2(part)-AUX1(have.sg)	no	yes	no	no	no	yes	no	no	yes
INF3-MOD2(inf)-AUX1(have.sg)	yes	yes	yes	no	yes	no	yes	yes	no

- ▶ starting point: a  $31 \times 267$  data table whereby each cluster order represents a row and each dialect location a column
- ▶ the dialect locations are now used to determine the degree of difference/similarity between the various cluster orders → each of the 31 cluster orders is compared to each other cluster order on 267 variables (i.e. as many as there are dialect locations)
- ▶ when we reduce the 31-dimensional distance matrix to a two-dimensional space, we can plot the differences and similarities between the cluster orders from the SAND-project

This talk in one  
slide

Introduction

The data: dialect  
Dutch verb clusters

Research questions

Theoretical  
background:  
dialectometry

Methodology:  
reverse  
dialectometry

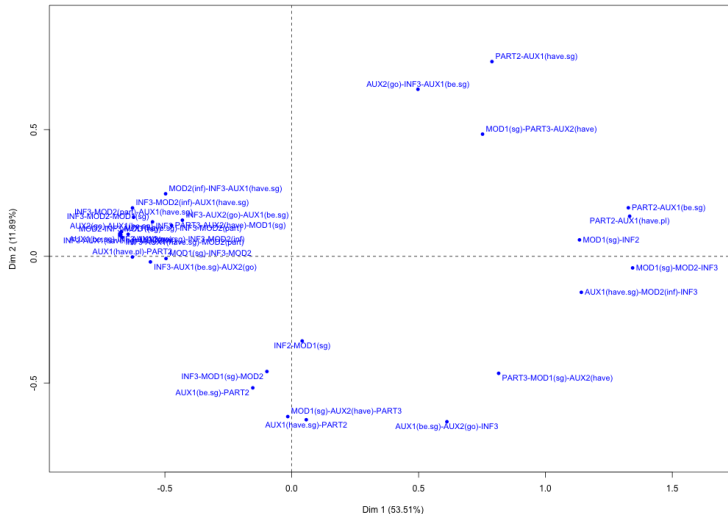
Results

Main conclusion

Extensions

References

### Two-dimensional representation of the 31 SAND-verb cluster orders



Methodology:  
reverse  
dialectometry

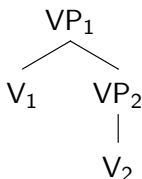
## Extensions

## References

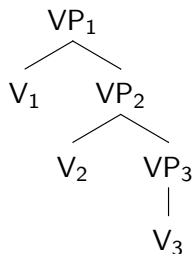
- ▶ **note:** each point now represents a particular cluster order and closeness of points indicates how alike two verb cluster orders are based on their geographical spread
- ▶ if this likeness is the result of grammar, i.e. grammatical microparameters, then verb cluster orders that are near one another should be the result of the same parameter setting, i.e. parameters create 'natural classes' of verb cluster orders
- ▶ in order to find those parameters, we can also encode the cluster orders in terms of their theoretical linguistic analyses
- ▶ example: Barbiers (2005)

- ▶ Barbiers (2005) derives verb cluster orders as follows:
  - ▶ base order is uniformly head-initial → derives 12 and 123

(10)



(11)



This talk in one  
slide

Introduction

The data: dialect  
Dutch verb clusters

Research questions

Theoretical  
background:  
dialectometry

Methodology:  
reverse  
dialectometry

Results

Main conclusion

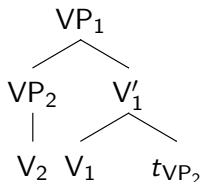
Extensions

References

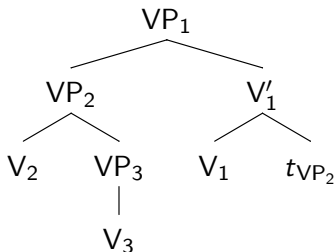


- ▶ Barbiers (2005) derives verb cluster orders as follows:
  - ▶ movement is VP-intrapolation → derives 21 and 231, 312 and 132, and fails to derive 213

(12)



(13)



This talk in one  
slide

Introduction

The data: dialect  
Dutch verb clusters

Research questions

Theoretical  
background:  
dialectometry

Methodology:  
reverse  
dialectometry

Results

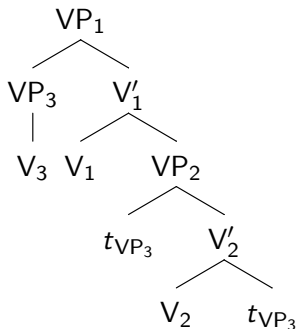
Main conclusion

Extensions

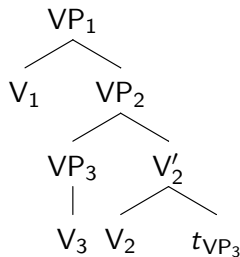
References

- Barbiers (2005) derives verb cluster orders as follows:
  - movement is VP-intrapolation → derives 21 and 231, 312 and 132, and fails to derive 213

(14)



(15)



This talk in one  
slide

Introduction

The data: dialect  
Dutch verb clusters

Research questions

Theoretical  
background:  
dialectometry

Methodology:  
reverse  
dialectometry

Results

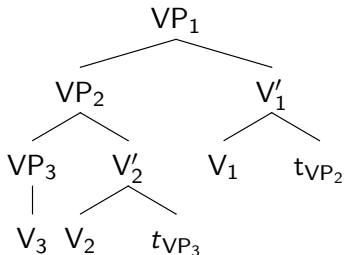
Main conclusion

Extensions

References

- ▶ Barbiers (2005) derives verb cluster orders as follows:
  - ▶ VP-intrapolation can pied-pipe other material → derives 321 (movement of VP3 to specVP1 via specVP2 and with pied-piping of VP2)

(16)



- ▶ Barbiers (2005) derives verb cluster orders as follows:
  - ▶ VP intraposition is triggered by feature checking: modal and aspectual auxiliaries enter into a(n eventive) feature checking relation with the main verb, while perfective auxiliaries enter into a perfective checking relationship with their immediately selected verb → rules out 231 in the case of MOD-MOD/AUX-V-clusters and 312 in the case of AUX-AUX/MOD-V-clusters

(17)  $[VP_1 \text{ mod}_{[uEvent]} [VP_2 \text{ mod}_{[uEvent]} [VP_3 \text{ inf}_{[iEvent]} ]]]$

(18)  $[VP_1 \text{ aux}_{[uPerf]} [VP_2 \text{ mod}_{[iPerf, uEvent]} [VP_3 \text{ inf}_{[iEvent]} ]]]$

- ▶ from this theoretical account we can distill the following micro-parameters:
  - ▶ [ $\pm$ base-generation]: can the order be base-generated?
  - ▶ [ $\pm$ movement]: can the order be derived via movement?
  - ▶ [ $\pm$ pied-piping]: does the derivation involve pied-piping?
  - ▶ [ $\pm$ feature-checking violation]: does the order involve a feature checking violation?
- ▶ and the 31 SAND cluster orders can be encoded in terms of these micro-parameters

This talk in one  
slide

Introduction

The data: dialect  
Dutch verb clusters

Research questions

Theoretical  
background:  
dialectometry

Methodology:  
reverse  
dialectometry

Results

Main conclusion

Extensions

References

	Barbiers-base.generation	Barbiers-movement	Barbiers-spec-pied-piping	Barbiers-feature.checking-failure
AUX1(be.sg)-PART2	yesBase	noMvt	noPiedP	noFeatCheckFail
PART2-AUX1(be.sg)	noBase	yesMvt	noPiedP	noFeatCheckFail
AUX1(have.sg)-PART2	yesBase	noMvt	noPiedP	noFeatCheckFail
PART2-AUX1(have.sg)	noBase	yesMvt	noPiedP	noFeatCheckFail
AUX1(have.pl)-PART2	yesBase	noMvt	noPiedP	noFeatCheckFail
PART2-AUX1(have.pl)	noBase	yesMvt	noPiedP	noFeatCheckFail
MOD1(sg)-INF2	yesBase	noMvt	noPiedP	noFeatCheckFail
INF2-MOD1(sg)	noBase	yesMvt	noPiedP	noFeatCheckFail
MOD2-INF3-MOD1(sg)	noBase	yesMvt	noPiedP	yesFeatCheckFail
MOD1(sg)-MOD2-INF3	yesBase	noMvt	noPiedP	noFeatCheckFail
MOD1(sg)-INF3-MOD2	noBase	yesMvt	noPiedP	noFeatCheckFail
INF3-MOD2-MOD1(sg)	noBase	yesMvt	yesPiedP	noFeatCheckFail
INF3-MOD1(sg)-MOD2	noBase	yesMvt	noPiedP	noFeatCheckFail
MOD1(sg)-AUX2(have)-PART3	yesBase	noMvt	noPiedP	noFeatCheckFail
MOD1(sg)-PART3-AUX2(have)	noBase	yesMvt	noPiedP	noFeatCheckFail
PART3-MOD1(sg)-AUX2(have)	noBase	yesMvt	noPiedP	noFeatCheckFail
PART3-AUX2(have)-MOD1(sg)	noBase	yesMvt	yesPiedP	noFeatCheckFail
AUX1(be.sg)-AUX2(go)-INF3	yesBase	noMvt	noPiedP	noFeatCheckFail
AUX1(be.sg)-INF3-AUX2(go)	noBase	yesMvt	noPiedP	noFeatCheckFail
AUX2(go)-AUX1(be.sg)-INF3	noBase	noMvt	noPiedP	noFeatCheckFail
AUX2(go)-INF3-AUX1(be.sg)	noBase	yesMvt	noPiedP	noFeatCheckFail
INF3-AUX1(be.sg)-AUX2(go)	noBase	yesMvt	noPiedP	yesFeatCheckFail
INF3-AUX2(go)-AUX1(be.sg)	noBase	yesMvt	noPiedP	noFeatCheckFail
AUX1(have.sg)-MOD2(inf)-INF3	yesBase	noMvt	noPiedP	noFeatCheckFail
AUX1(have.sg)-INF3-MOD2(part)	noBase	yesMvt	noPiedP	noFeatCheckFail
AUX1(have.sg)-INF3-MOD2(inf)	noBase	yesMvt	noPiedP	noFeatCheckFail
MOD2(inf)-INF3-AUX1(have.sg)	noBase	yesMvt	noPiedP	noFeatCheckFail
INF3-AUX1(have.sg)-MOD2(inf)	noBase	yesMvt	noPiedP	yesFeatCheckFail
INF3-AUX1(have.sg)-MOD2(part)	noBase	yesMvt	noPiedP	yesFeatCheckFail
INF3-MOD2(part)-AUX1(have.sg)	noBase	yesMvt	noPiedP	noFeatCheckFail
INF3-MOD2(inf)-AUX1(have.sg)	noBase	yesMvt	noPiedP	noFeatCheckFail

Quantity and  
quality in linguistic  
variation

Jeroen van  
Craenenbroeck

This talk in one  
slide

Introduction

The data: dialect  
Dutch verb clusters

Research questions

Theoretical  
background:  
dialectometry

Methodology:  
reverse  
dialectometry

Results

Main conclusion

Extensions

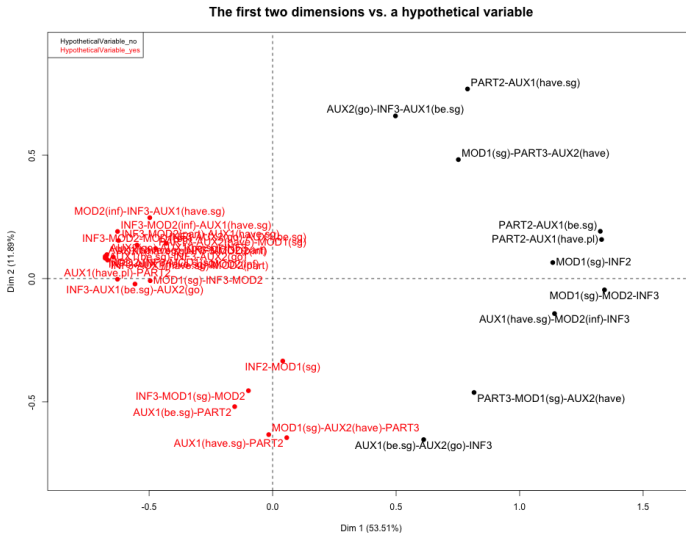
References

- ▶ Barbiers's microparameters thus serve as additional variables in our data table
- ▶ **in total:** 70 additional variables distilled from the theoretical literature on verb clusters:
  - ▶ the analyses of Barbiers (2005), Barbiers and Bennis (2010), Abels (2011), Haegeman and Riemsdijk (1986), Bader (2012), and Schmid and Vogel (2004)
  - ▶ a head-initial head movement analysis, a head-final head movement analysis, a head-initial XP-movement analysis, a head-final XP-movement analysis (all based on Wurmbrand (2005))
  - ▶ 17 additional variables based on the theoretical literature, but not linked to a specific analysis
- ▶ **note:** all these variables are encoded as **supplementary** variables, i.e. they do not contribute to the construction of the verb cluster plot, but they can be mapped/matched against it

- ▶ there are three ways of testing how well a supplementary (linguistic) variable lines up with the output of the geographical analysis:

1. visual inspection of a color-coded plot





- ▶ there are three ways of testing how well a supplementary (linguistic) variable lines up with the output of the geographical analysis:

1. visual inspection of a color-coded plot
2.  $\eta^2$  (squared correlation ratio): value between 0 and 1 indicating the strength of the link between the dimension and a particular categorical variable; can be interpreted as the percentage of variation on the dimension that can be explained by the categorical variable

	dimension 1	dimension 2
HypotheticalVariable	0.861	0.043

- ▶ **word of caution:**  $\eta^2$  also goes up if the number of possible values of the categorical variable goes up (Richardson (2011)) → safest option is to look for variables with a high  $\eta^2$  **and** only two or three possible values

- there are three ways of testing how well a supplementary (linguistic) variable lines up with the output of the geographical analysis:

1. visual inspection of a color-coded plot
2.  $\eta^2$  (squared correlation ratio): value between 0 and 1 indicating the strength of the link between the dimension and a particular categorical variable; can be interpreted as the percentage of variation on the dimension that can be explained by the categorical variable
3. v-test: value indicating whether a value of a categorical variable varies significantly from 0 on a particular dimension; significant values are  $< -2$  or  $> 2$

	<b>dimension 1</b>	<b>dimension 2</b>
HypotheticalVariable_yes	-5.082	-1.130
HypotheticalVariable_no	5.082	1.130

# Results

- ▶ **recall:** we are trying to determine if the variation in word order in verbal clusters is determined by grammatical parameters, and if so to what extent
- ▶ this means we need to determine **how many** parameters there are and **what they are**
- ▶ **proposal (I):** the number of parameters responsible for the verb cluster variation = the number of dimensions we reduce our data set to
- ▶ **proposal (II):** the identity of those parameters = the interpretation of the dimensions

Quantity and  
quality in linguistic  
variation

Jeroen van  
Craenenbroeck

This talk in one  
slide

Introduction

The data: dialect  
Dutch verb clusters

Research questions

Theoretical  
background:  
dialectometry

Methodology:  
reverse  
dialectometry

Results

Introduction

The number of  
relevant dimensions

Dimension 1

Dimension 2

Dimension 3

Conclusion

Main conclusion

Extensions

References

# The number of relevant dimensions

- ▶ recall: the analysis reduces a multi-dimensional distance matrix into a low-dimensional space, *while retaining as much as possible of the information present in the original object*
- ▶ we can plot the percentage of variance explained per dimension (= scree plot)

Quantity and  
quality in linguistic  
variation

Jeroen van  
Craenenbroeck

This talk in one  
slide

Introduction

The data: dialect  
Dutch verb clusters

Research questions

Theoretical  
background:  
dialectometry

Methodology:  
reverse  
dialectometry

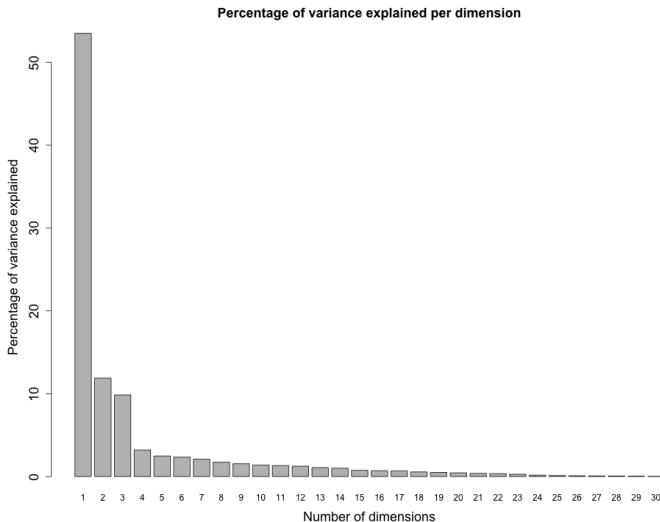
Results

Introduction  
**The number of  
relevant dimensions**  
Dimension 1  
Dimension 2  
Dimension 3  
Conclusion

Main conclusion

Extensions

References



- ▶ **note:** there seems to be a clear cut-off point after the third dimension
- ▶ together, the first three dimensions account for 78.46% of the variation in the SAND verb cluster data
- ▶ this means that roughly 80% of the variation in verb cluster ordering in SAND can be reduced to three parameters
- ▶ in order to know what those parameters are, we need to interpret the first three dimensions

This talk in one  
slide

Introduction

The data: dialect  
Dutch verb clusters

Research questions

Theoretical  
background:  
dialectometry

Methodology:  
reverse  
dialectometry

Results

Introduction

The number of  
relevant dimensions

Dimension 1

Dimension 2

Dimension 3

Conclusion

Main conclusion

Extensions

References

# Dimension 1

- ▶ highest  $\eta^2$ -values:

## dimension 1

BarBen.NomInf	0.425
Bader.VMod	0.398

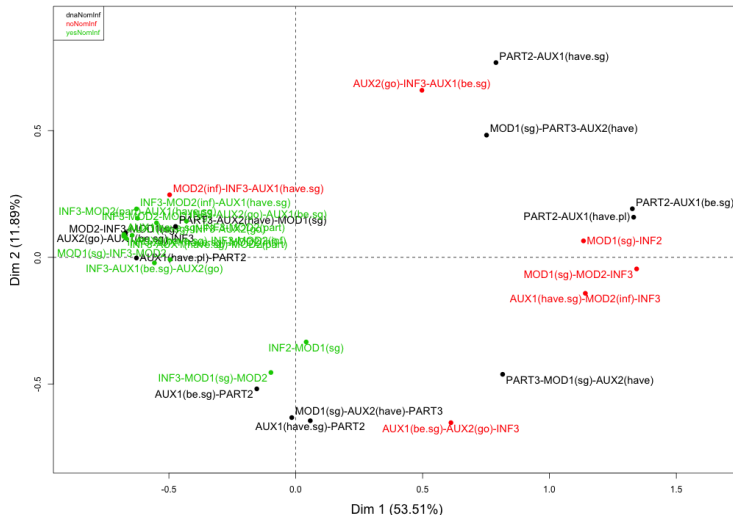
- ▶ BarBen.NomInf: Barbiers and Bennis (2010): the infinitival main verb is nominalized
  - ▶ Bader.VMod: Bader (2012): the complement of a modal verb precedes the modal
- ▶ confirmed by v-test:

## v-test dimension 1

yesNomInf	-3.303
yesVMod	-3.090



Dimension 1 vs. Barbiers & Bennis's (2010) nominalized infinitives



This talk in one  
slide

Introduction

The data: dialect  
Dutch verb clusters

Research questions

Theoretical  
background:  
dialectometry

Methodology:  
reverse  
dialectometry

Results

Introduction  
The number of  
relevant dimensions

**Dimension 1**

Dimension 2

Dimension 3

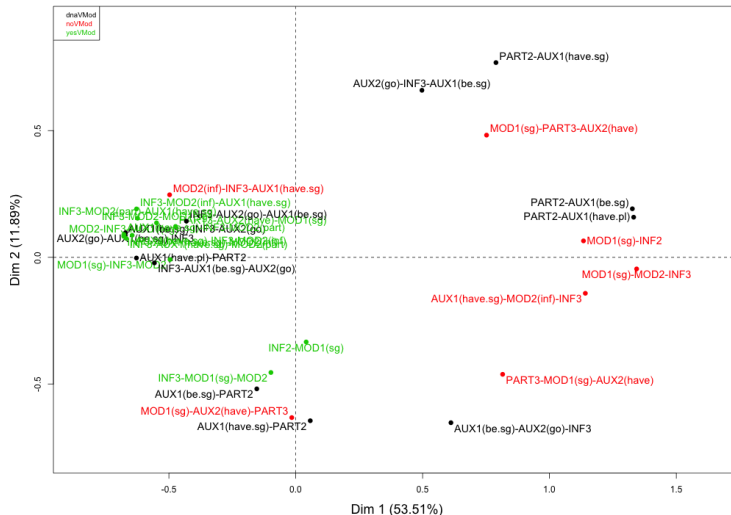
Conclusion

Main conclusion

Extensions

References

Dimension 1 vs. Bader's (2012) VMod-constraint



This talk in one  
slide

Introduction

The data: dialect  
Dutch verb clusters

Research questions

Theoretical  
background:  
dialectometry

Methodology:  
reverse  
dialectometry

Results

Introduction  
The number of  
relevant dimensions

**Dimension 1**

Dimension 2

Dimension 3

Conclusion

Main conclusion

Extensions

References

- ▶ **note:** while both variables propose a general split that seems well represented on dimension 1, there are a number of verb clusters orders for which they are irrelevant (because the cluster doesn't contain the relevant configuration)
- ▶ let's try to strengthen our interpretation of dimension 1 by incorporating those points
- ▶ Abels (2011): there seems to be a correlation between the position of infinitives *vis-à-vis* modals on the one hand and the position of participles *vis-à-vis* auxiliaries on the other

This talk in one  
slide

Introduction

The data: dialect  
Dutch verb clusters

Research questions

Theoretical  
background:  
dialectometry

Methodology:  
reverse  
dialectometry

Results

Introduction  
The number of  
relevant dimensions  
**Dimension 1**  
Dimension 2  
Dimension 3  
Conclusion

Main conclusion

Extensions

References

- ▶ new variable: InfMod.AuxPart:
  - ▶ set to 'no' when the modal precedes the infinitive (when present) and the participle precedes the auxiliary (when present)
  - ▶ set to 'yes' when at least one of these conditions is not met
- ▶  $\eta^2$  of InfMod.AuxPart: 0.6142

v-test dimension 1	
InfMod.AuxPart_yes	-4.292
InfMod.AuxPart_no	4.292

This talk in one  
slide

Introduction

The data: dialect  
Dutch verb clusters

Research questions

Theoretical  
background:  
dialectometryMethodology:  
reverse  
dialectometry

Results

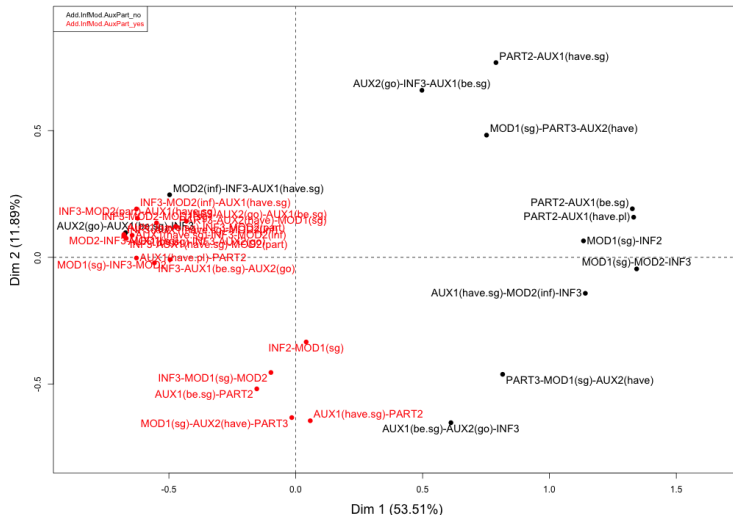
Introduction  
The number of  
relevant dimensions  
**Dimension 1**  
Dimension 2  
Dimension 3  
Conclusion

Main conclusion

Extensions

References

Dimension 1 vs. the new InfMod.AuxPart-variable



This talk in one  
slide

Introduction

The data: dialect  
Dutch verb clusters

Research questions

Theoretical  
background:  
dialectometry

Methodology:  
reverse  
dialectometry

Results

Introduction  
The number of  
relevant dimensions

**Dimension 1**

Dimension 2

Dimension 3

Conclusion

Main conclusion

Extensions

References

- ▶ **note:** this new variable aligns very nicely with the first dimension
  - ▶ two recalcitrant cluster orders:
    - ▶ AUX2(go)-AUX1(be.sg)-INF3: possibly spurious: this is an order that seems excluded in any cluster in Dutch (only two hits in the whole of SAND)
    - ▶ MOD2(inf)-INF3-AUX1(have.sg)
- ▶ **this means** that the first (and most important) source of variation in Dutch verb clusters—i.e. the first microparameter—concerns the placement of modals and auxiliaries vs. the verbs they select
- ▶ it sets apart dialects that consistently place infinitives to the right and participles to the left from those that don't

## Dimension 2

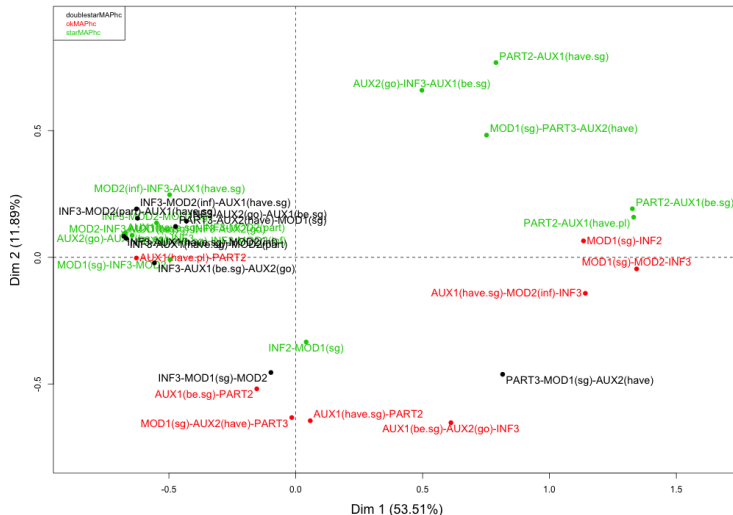
- ▶ highest  $\eta^2$ -values:

<b>dimension 2</b>	
SchmiVo.MAPhc	0.379
Barbiers.base.generation	0.309

- ▶ SchmiVo.MAPhc: Schmid and Vogel (2004): “If A and B are sister nodes at LF, and A is a head and B is a complement, then the correspondent of A precedes the one of B at PF.”
  - ▶ Barbiers.base.generation: Barbiers (2005): head-initial base structure
- ▶ confirmed by v-test:

<b>v-test dimension 2</b>	
Barbiers.base.generation_no	3.044
SchmiVo.MAPhc_starMAPhc	2.855
SchmiVo.MAPhc_okMAPhc	-3.044

Dimension 2 vs. Schmid & Vogel's (2004) MAPhc-constraint



This talk in one  
slide

Introduction

The data: dialect  
Dutch verb clusters

Research questions

Theoretical  
background:  
dialectometry

Methodology:  
reverse  
dialectometry

Results

Introduction  
The number of  
relevant dimensions  
Dimension 1  
**Dimension 2**  
Dimension 3  
Conclusion

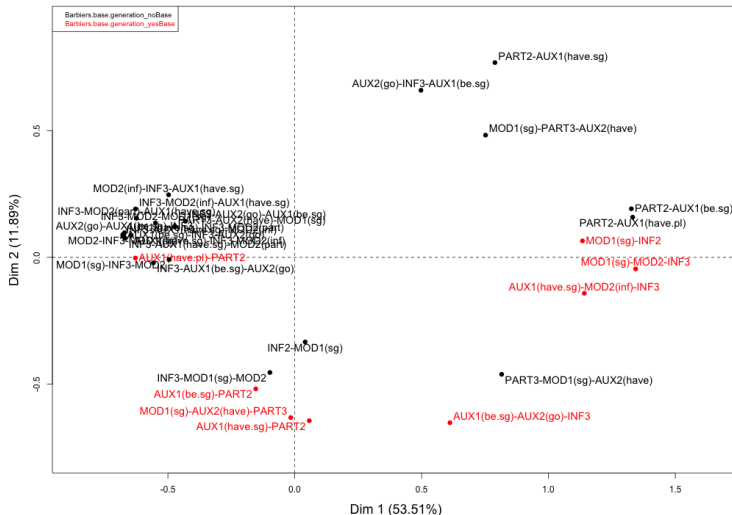
Main conclusion

Extensions

References



Dimensions 1 and 2 of the verb cluster MCA vs. Barbiers's (2005) base-generation



This talk in one  
slide

Introduction

The data: dialect  
Dutch verb clusters

Research questions

Theoretical  
background:  
dialectometry

Methodology:  
reverse  
dialectometry

Results

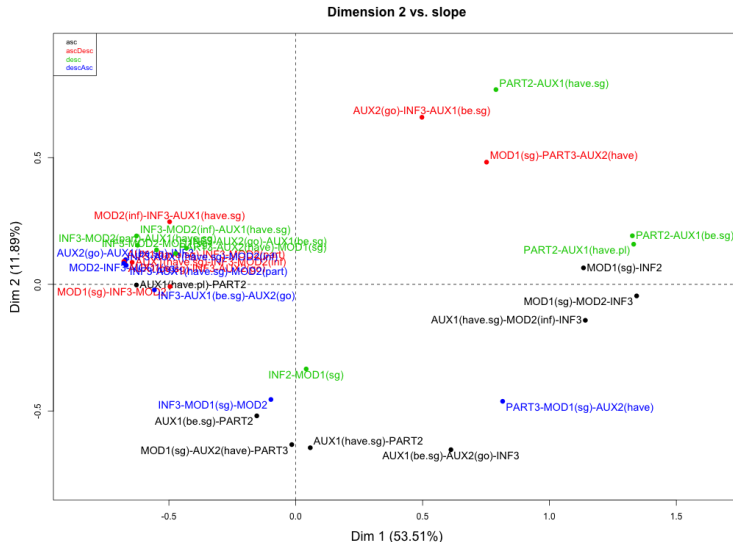
Introduction  
The number of  
relevant dimensions  
Dimension 1  
**Dimension 2**  
Dimension 3  
Conclusion

Main conclusion

Extensions

References

- ▶ **note:** just as was the case with dimension 1, the variables culled from the literature leave room for improvement in interpreting dimension 2
- ▶ another variable that does well is slope ( $\eta^2 = 0.422$ ): is the order ascending, descending, first-ascending-then-descending, or first-descending-then-ascending?



- ▶ **note:** ascDesc and desc pattern towards the positive values of dimension 2, while asc and descAsc tend to yield negative values for this dimension
- ▶ new variable: FinalDescent:
  - ▶ set to 'yes' if the cluster ends in a descending order
  - ▶ set to 'no' if it ends in an ascending order

FinalDescent_yes	FinalDescent_no
21	12
132	123
321	312
231	213

This talk in one  
slide

Introduction

The data: dialect  
Dutch verb clusters

Research questions

Theoretical  
background:  
dialectometryMethodology:  
reverse  
dialectometry

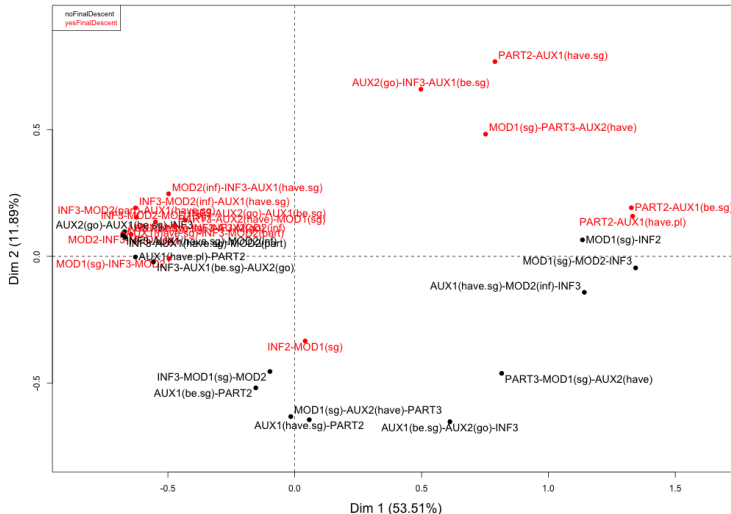
Results

Introduction  
The number of  
relevant dimensions  
Dimension 1  
**Dimension 2**  
Dimension 3  
Conclusion

Main conclusion

Extensions

References



- ▶  $\eta^2$  of FinalDescent: 0.382

**v-test dimension 2**

FinalDescent_yes	3.387
FinalDescent_no	-3.387

- ▶ possible theoretical interpretation of FinalDescent: it groups together cluster orders which are 0 or 1 movement operations away from a strictly head-final order (i.e. 132, 321, 231), from those that require at least two movement operations (123, 312, 213)
  - ▶ caveat: two-verb clusters → there are only two possible orders, so you can always get from one to the other with one movement operation
- ▶ **this means** that the second source of variation in Dutch verb clusters—i.e. the second microparameter—concerns the degree to which a cluster order diverges from a strictly head-final order

[This talk in one slide](#)[Introduction](#)[The data: dialect Dutch verb clusters](#)[Research questions](#)[Theoretical background: dialectometry](#)[Methodology: reverse dialectometry](#)[Results](#)[Introduction](#)  
[The number of relevant dimensions](#)  
[Dimension 1](#)  
[Dimension 2](#)  
[Dimension 3](#)  
[Conclusion](#)[Main conclusion](#)[Extensions](#)[References](#)

# Dimension 3

- ▶ highest  $\eta^2$ -values:

## dimension 3

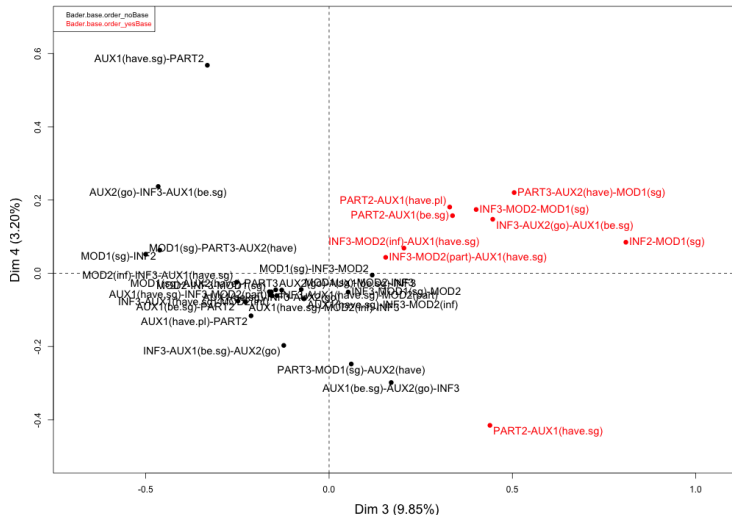
SchmiVo.MAPch	0.701
Bader.base.order	0.686

- ▶ SchmiVo.MAPch: Schmid and Vogel (2004): “If A and B are sister nodes at LF, and A is a head and B is a complement, then the correspondent of B precedes the one of A at PF.”
  - ▶ Bader.base.order: Bader (2012): a strictly head-final base order
- ▶ confirmed by v-test:

## v-test dimension 3

SchmiVo.MAPch_okMAPch	4.537
Bader.base.order_yes	4.537

Dimension 3 vs. Bader's (2012) base-generated order



This talk in one  
slide

Introduction

The data: dialect  
Dutch verb clusters

Research questions

Theoretical  
background:  
dialectometry

Methodology:  
reverse  
dialectometry

Results

Introduction  
The number of  
relevant dimensions  
Dimension 1  
Dimension 2  
**Dimension 3**  
Conclusion

Main conclusion

Extensions

References



- ▶ there is a very strong correlation between a head-final base order and the third dimension in the analysis
- ▶ **this means** that the third source of variation in Dutch verb clusters—i.e. the third microparameter—concerns the question of whether a dialect diverges from a strictly head final order or not

# Results: Conclusion

- ▶ interpreting the first three dimensions of the quantitative analysis of the verb cluster data in the Syntactic Atlas of Dutch Dialects allows us to construct the following (rough) parametric account of verb cluster ordering:
  1. a head-final base order
  2. which dialects can diverge from or not: [ $\pm$ Movement] (dimension 3)
  3. those that diverge can diverge strongly or not: Economy of Movement (dimension 2)
  4. above and beyond all this, a headedness parameter regulates the order of infinitives and participles *vis-à-vis* their selecting verbs: [ $\pm$ ModInf&PartAux] (dimension 1)

Quantity and  
quality in linguistic  
variation

Jeroen van  
Craenenbroeck

This talk in one  
slide

Introduction

The data: dialect  
Dutch verb clusters

Research questions

Theoretical  
background:  
dialectometry

Methodology:  
reverse  
dialectometry

Results

Introduction  
The number of  
relevant dimensions  
Dimension 1  
Dimension 2  
Dimension 3  
Conclusion

Main conclusion

Extensions

References

# Main conclusion

- ▶ the considerable variation found in Dutch verb cluster orders can be reduced to three grammatical microparameters:
  1. the order of modals and auxiliaries vs. the verbs they select:  $[\pm\text{ModInf}\&\text{PartAux}]$
  2. the degree of divergence from a head-final order:  $[\text{EconomyOfMovement}]$
  3. adherence to a head-final order or not:  $[\pm\text{Movement}]$
- ▶ more generally, there is room for fruitful collaboration between formal-theoretical and quantitative-statistical linguistics:
  - ▶ the former can guide the interpretation of results from the latter
  - ▶ the latter can help evaluate and test hypotheses of the former

# Extensions

## ► empirical (I):

- extend the analysis with more data from Dutch (written SAND-questionnaire) and other Germanic languages and varieties such as German and Swiss German

## ► empirical (II):

- apply these same techniques to other empirical domains

## ► theoretical:

- combine the three microparameters uncovered by the quantitative analysis into one, coherent theoretical account of verb clusters
- use additional statistical techniques (cluster analysis, association rule mining) to further explore the data
- find a way to compare not just individual theoretical hypotheses but rather entire analyses of verb cluster ordering

Quantity and  
quality in linguistic  
variation

Jeroen van  
Craenenbroeck

This talk in one  
slide

Introduction

The data: dialect  
Dutch verb clusters

Research questions

Theoretical  
background:  
dialectometry

Methodology:  
reverse  
dialectometry

Results

Main conclusion

Extensions

References

# References I

- Abels, Klaus. 2011. Hierarchy-order relations in the germanic verb cluster and in the noun phrase. *GAGL* 53:1–28.
- Bader, Markus. 2012. Verb-cluster variations: a harmonic grammar analysis. Handout of a talk presented at “New ways of analyzing syntactic variation”, November 2012.
- Barbiers, Sjef. 2005. Word order variation in three-verb clusters and the division of labour between generative linguistics and sociolinguistics. In *Syntax and variation. Reconciling the biological and the social*, ed. Leonie Cornips and Karen P. Corrigan, volume 265 of *Current issues in linguistic theory*, 233–264. John Benjamins.
- Barbiers, Sjef, and Hans Bennis. 2010. De plaats van het werkwoord in zuid en noord. In *Voor Magda. Artikelen voor Magda Devos bij haar afscheid van de Universiteit Gent*, ed. Johan De Caluwe and Jacques Van Keymeulen, 25–42. Gent: Academia.
- Borg, Ingwer, and Patrick J.F. Groenen. 2005. *Modern Multidimensional Scaling. Theory and applications.*. Springer, 2nd edition.
- Haegeman, Liliane, and Henk van Riemsdijk. 1986. Verb projection raising, scope, and the typology of verb movement rules. *Linguistic Inquiry* 17:417–466.

Quantity and  
quality in linguistic  
variation

Jeroen van  
Craenenbroeck

This talk in one  
slide

Introduction

The data: dialect  
Dutch verb clusters

Research questions

Theoretical  
background:  
dialectometry

Methodology:  
reverse  
dialectometry

Results

Main conclusion

Extensions

References

# References II

- Nerbonne, John, and William A. Kretzschmar Jr. 2013. Dialectometry++. *Literary and Linguistic Computing* 28:2–12.
- Richardson, John T.E. 2011. Eta squared and partial eta squared as measures of effect size in educational research. *Educational Research Review* 6:135–147.
- Schmid, Tanja, and Ralf Vogel. 2004. Dialectal variation in German 3-Verb clusters. *The Journal of Comparative Germanic Linguistics* 7:235–274.
- Wurmbrand, Susanne. 2005. Verb clusters, verb raising, and restructuring. In *The Blackwell Companion to Syntax*, ed. Martin Everaert and Henk van Riemsdijk, volume V, chapter 75, 227–341. Oxford: Blackwell.

Quantity and  
quality in linguistic  
variation

Jeroen van  
Craenenbroeck

This talk in one  
slide

Introduction

The data: dialect  
Dutch verb clusters

Research questions

Theoretical  
background:  
dialectometry

Methodology:  
reverse  
dialectometry

Results

Main conclusion

Extensions

References