

Mapping language varieties

Koen Plevoets

Faculty of Applied Language Studies

University College Ghent

Mapping language varieties

- The use of correspondence analysis to represent linguistic distances

Mapping language varieties

- The use of correspondence analysis to represent linguistic distances
- The varieties of Belgian Dutch

Mapping language varieties

- Belgian Dutch: a lot of linguistic variation
- Formal-informal

Mapping language varieties

- **OBJECTIVE:** Measure difference between ‘varieties’/‘registers’ of Belgian Dutch based on their degree of (in)formality

	Adm	Ext	Fic	Ins	Jour	Non_f
verkrijgen	187	144	1	33	28	17
bekomen	91	69	0	23	3	6
zodra	47	48	10	17	85	26
van.zodra	7	4	0	3	5	0
...						

Mapping language varieties

2 datasets:

- TSS: Standard vs. vernacular
 - Spoken Dutch Corpus (CGN)
 - 14 linguistic variables
 - Is ‘tussentaal’ the new omni-situational standard?
- COMURE: Translated vs. non-translated
 - Dutch Parallel Corpus
 - 13 linguistic variables
 - Is translated language more formal than non-translated language?

Overview

- Correspondence analysis: Chi-square as a distance metric
- Within-chi-square
- ‘Multifactorial’ correspondence analysis
- Statistical inference
- Conclusions

Correspondence analysis

- Frequency Table: Chi-square
- X^2 = measure of heterogeneity

$$X^2 = \sum_{i=1}^I \sum_{j=1}^J \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

$$E_{ij} = \frac{N_{i+} * N_{+j}}{N_{++}}$$

Correspondence analysis

$$E_{ij} = \frac{N_{i+} * N_{+j}}{N_{++}}$$

	Adm	Ext	Fic	Ins	Jour	Non_f	Ni+
verkrijgen	187	144	1	33	28	17	410
bekomen	91	69	0	23	3	6	192
zodra	47	48	10	17	85	26	233
van.zodra	7	4	0	3	5	0	19
N+j	332	265	11	76	121	49	854

Correspondence analysis

- Chi-square as a distance metric
- Anderson (2003: 80): X^2 of 2 varieties = Mahalanobis distance between the 2 varieties

$$D^2(Adm, Ext) = (Adm - Ext)^T * S^{-1} * (Adm - Ext)$$

- S = Covariance matrix of the linguistic items

Correspondence analysis

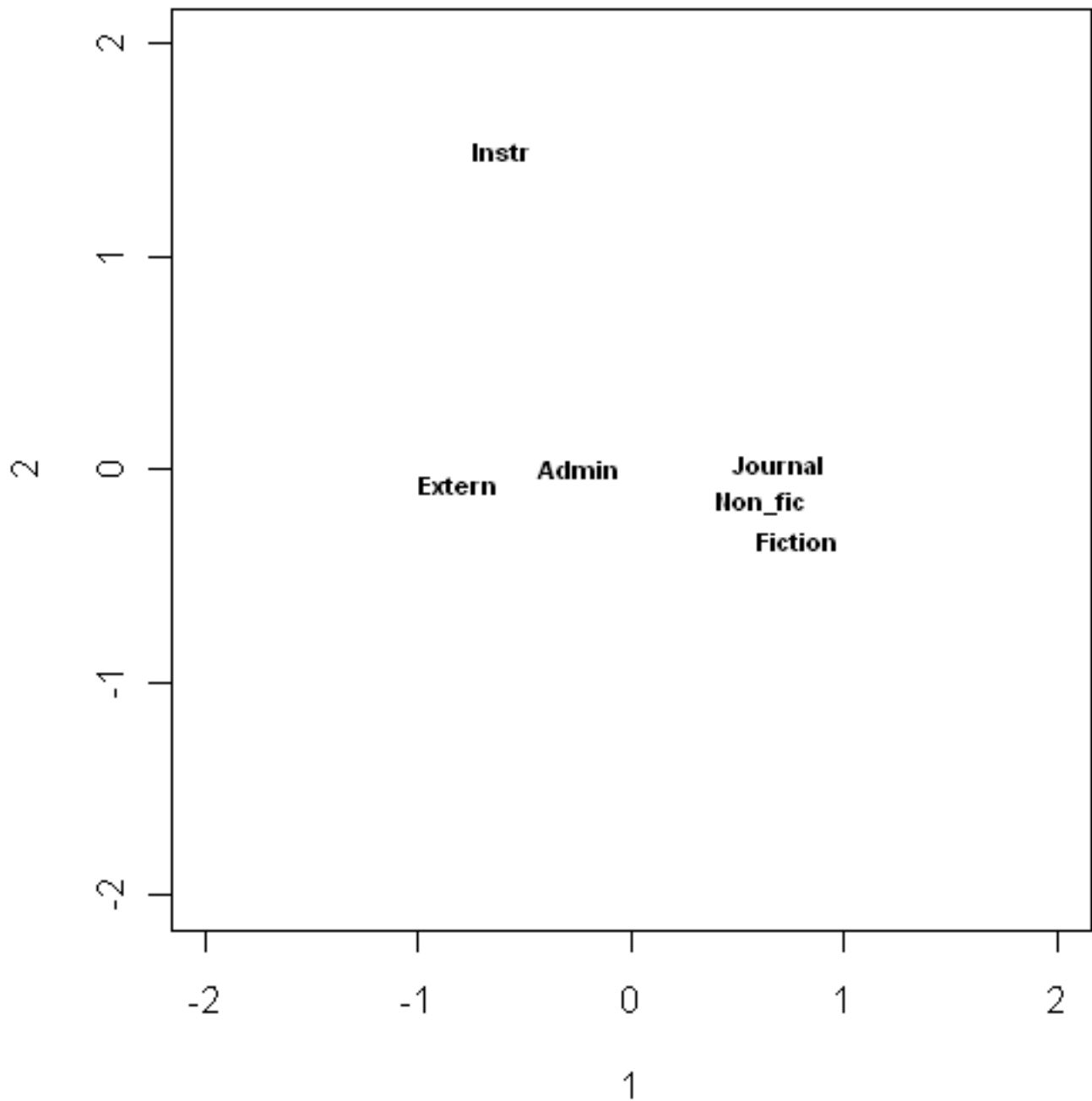
- OUTPUT: representation in low-dimensional space
- ‘Singular Value Decomposition’: $X = U * \Sigma * V^T$

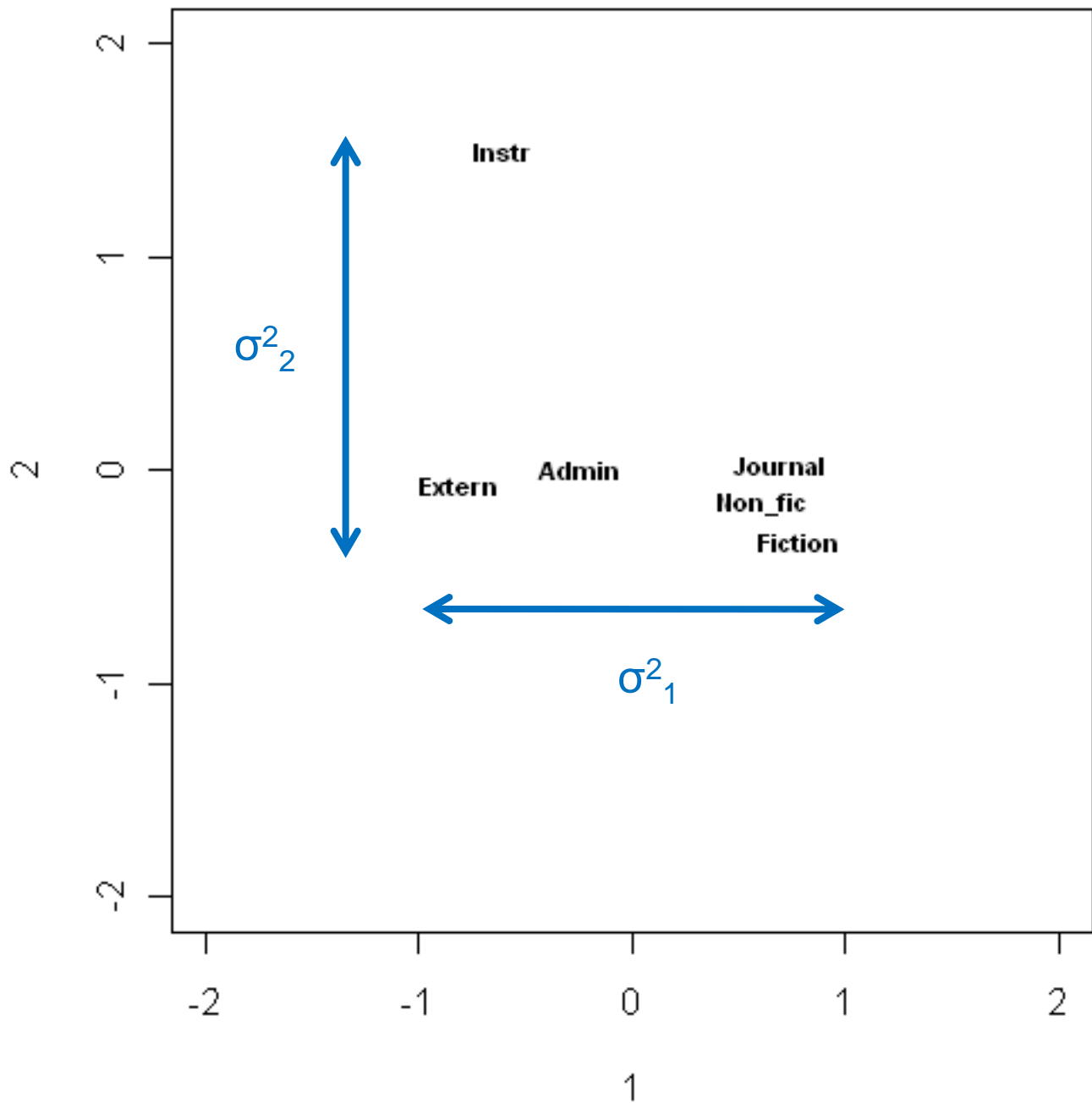
Correspondence analysis

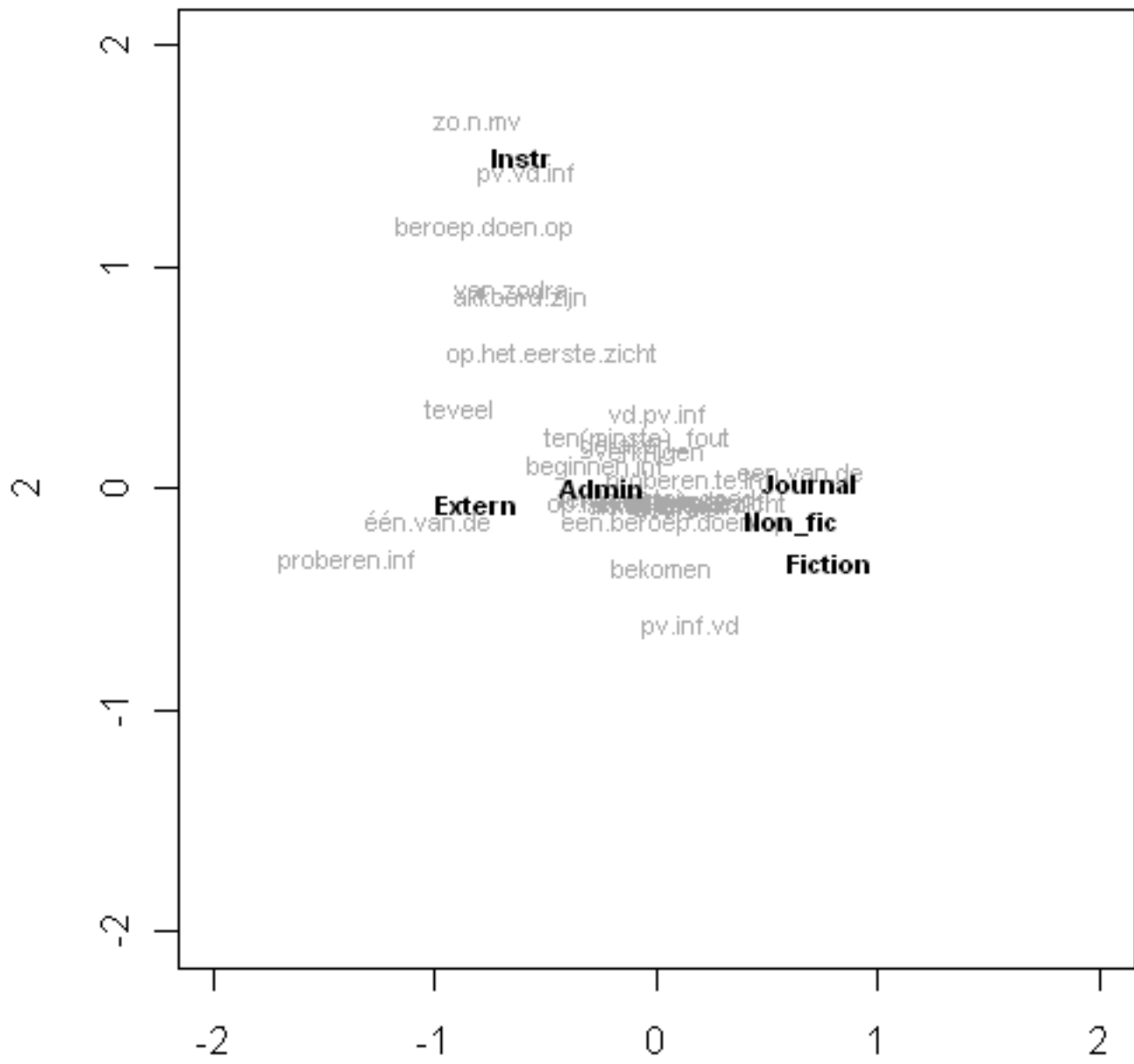
- OUTPUT: representation in low-dimensional space
- ‘Singular Value Decomposition’: $X = U * \Sigma * V^T$

$V =$

	1	2
Admin	-0.24385	0.00865
Extern	-0.80581	-0.05939
Fiction	0.79083	-0.32532
Instr	-0.59552	1.496837
Journal	0.70657	0.026645
Non_fic	0.611468	-0.14888

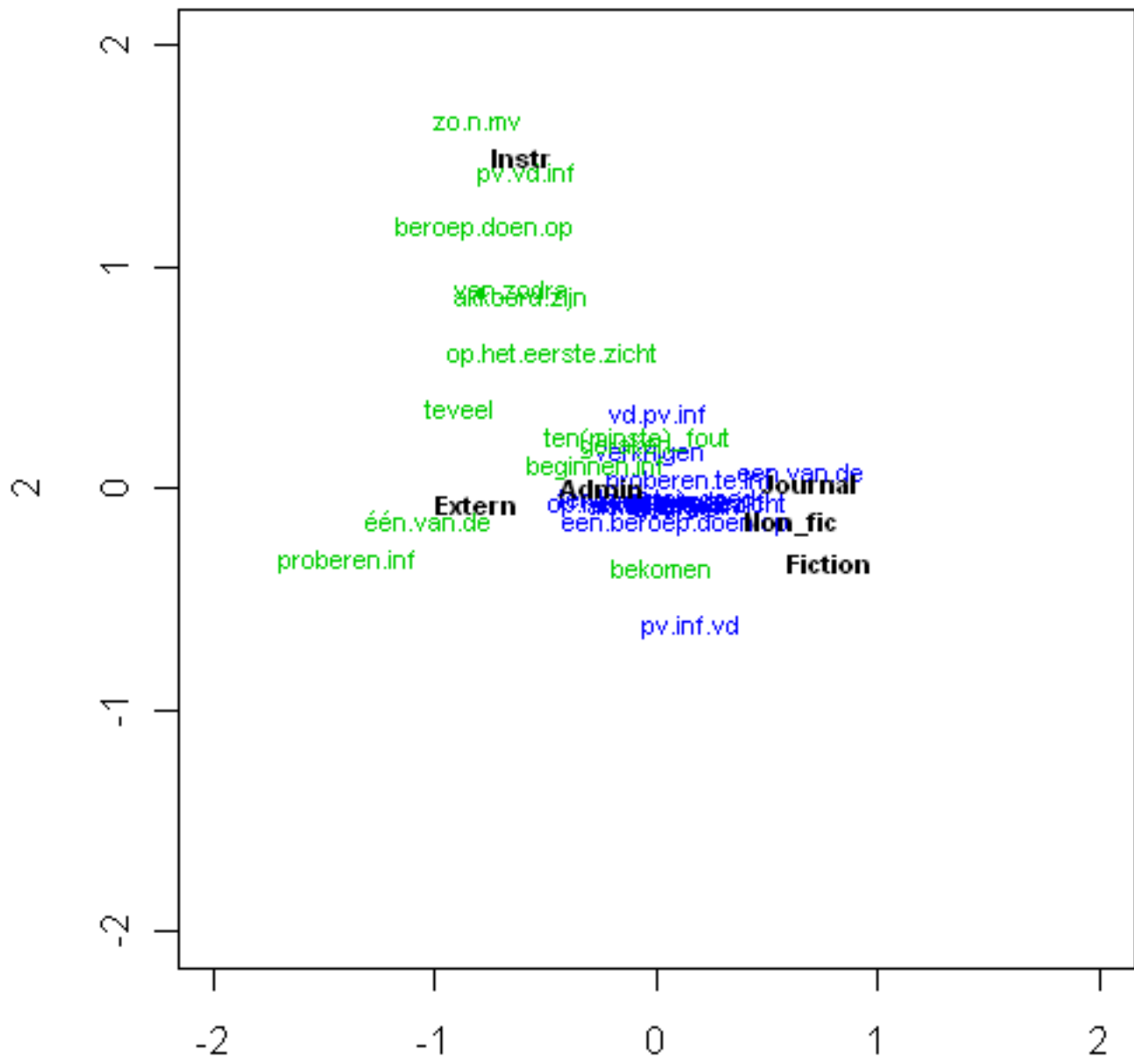






Biplot

- Gower, J., S. Lubbe & N. Le Roux (2011). *Understanding biplots*. Hoboken: Wiley.
- Greenacre, M. (2010). *Biplots in practice*. Bilbao: Fundación BBVA.



Within-chi-square

- ‘Linguistic variable’: alternation of linguistic variants
- HENCE: Linguistic items are partitioned

~ ‘Profile-based uniformity’ (Speelman, e.a. 2003)

Within-chi-square

	Adm	Ext	Fic	Ins	Jour	Non_f
verkrijgen	187	144	1	33	28	17
bekomen	91	69	0	23	3	6
zodra	47	48	10	17	85	26
van.zodra	7	4	0	3	5	0
pv.inf.vd	63	20	1	6	10	11
vd.pv.inf	20	14	1	9	19	1
pv.vd.inf	14	7	0	9	3	0
...						

Within-chi-square

	Adm	Ext	Fic	Ins	Jour	Non_f
verkrijgen	187	144	1	33	28	17
bekomen	91	69	0	23	3	6
zodra	47	48	10	17	85	26
van.zodra	7	4	0	3	5	0
pv.inf.vd	63	20	1	6	10	11
vd.pv.inf	20	14	1	9	19	1
pv.vd.inf	14	7	0	9	3	0
...						

Within-chi-square

- ‘Linguistic variable’: alternation of linguistic variants
- HENCE: Linguistic items are partitioned
- GOAL: Measure heterogeneity **per** variable (bracketing the heterogeneity across variables)
- ~ ‘Profile-based uniformity’ (Speelman, e.a. 2003)

Within-chi-square

- Huyghens' theorem (Greenacre 1984: 203-204)
- $X^2 = X^2_{\text{Within}} + X^2_{\text{Between}}$

$$X^2_{\text{Within}} = \sum_{k=1}^K \sum_{i=1}^{I_k} \sum_{j=1}^J \frac{\left(O_{ij} - \frac{N_{i+} * N_{+kj}}{N_{++}} \right)^2}{\frac{N_{i+} * N_{+j}}{N_{++}}}$$

$$X^2_{\text{Between}} = \sum_{k=1}^K \sum_{j=1}^J \frac{\left(N_{+kj} - \frac{N_{++} * N_{+j}}{N_{++}} \right)^2}{\frac{N_{++} * N_{+j}}{N_{++}}}$$

Within-chi-square

- Huyghens' theorem (Greenacre 1984: 203-204)
- $X^2 = X^2_{\text{Within}} + X^2_{\text{Between}}$

$$X^2_{\text{Within}} = \sum_{k=1}^K \sum_{i=1}^{I_k} \sum_{j=1}^J \frac{\left(O_{ij} - \frac{N_{i+} * N_{+k j}}{N_{++}} \right)^2}{\frac{N_{i+} * N_{+j}}{N_{++}}}$$

	Adm	Ext	Fic	Ins	Jour	Non_f	Ni+
verkrijgen	187	144	1	33	28	17	410
bekomen	91	69	0	23	3	6	192
zodra	47	48	10	17	85	26	233
van.zodra	7	4	0	3	5	0	19
...							

Within-chi-square

#	Variant 1	Variant 2	Variant 3	Translation
1	akkoord gaan	akkoord zijn		<i>to agree</i>
2	beginnen te +inf	beginnen + inf		<i>to begin to</i>
3	een beroep doen op	beroep doen op		<i>to appeal to</i>
4	een van de	één van de		<i>one of the</i>
5	op het eerste gezicht	op het eerste zicht		<i>at first sight</i>
6	proberen te + inf	proberen +inf		<i>to try to</i>
7	pv inf vd	vd pv inf	pv vd inf	<i>order of verbal end group</i>

Within-chi-square

#	Variant 1	Variant 2	Variant 3	Translation
8	raken	geraken		<i>to get</i>
9	te veel	teveel		<i>too much/many</i>
10	ten(minste)_goed	ten(minste)_fout		<i>at least</i>
11	verkrijgen	bekomen		<i>to obtain</i>
12	zodra	van zodra		<i>as soon as</i>
13	zulke + mv	zo'n + mv		<i>such + plural</i>

'Multifactorial' CA

- More than one factor coding the varieties
- Possibly with interactions

	DU_or	DU<EN	DU<FR
Adm	406	225	354
Ext	352	309	312
Fic	0	0	80
Ins	141	0	60
Jour	450	387	230
Non_f	363	0	93

'Multifactorial' CA

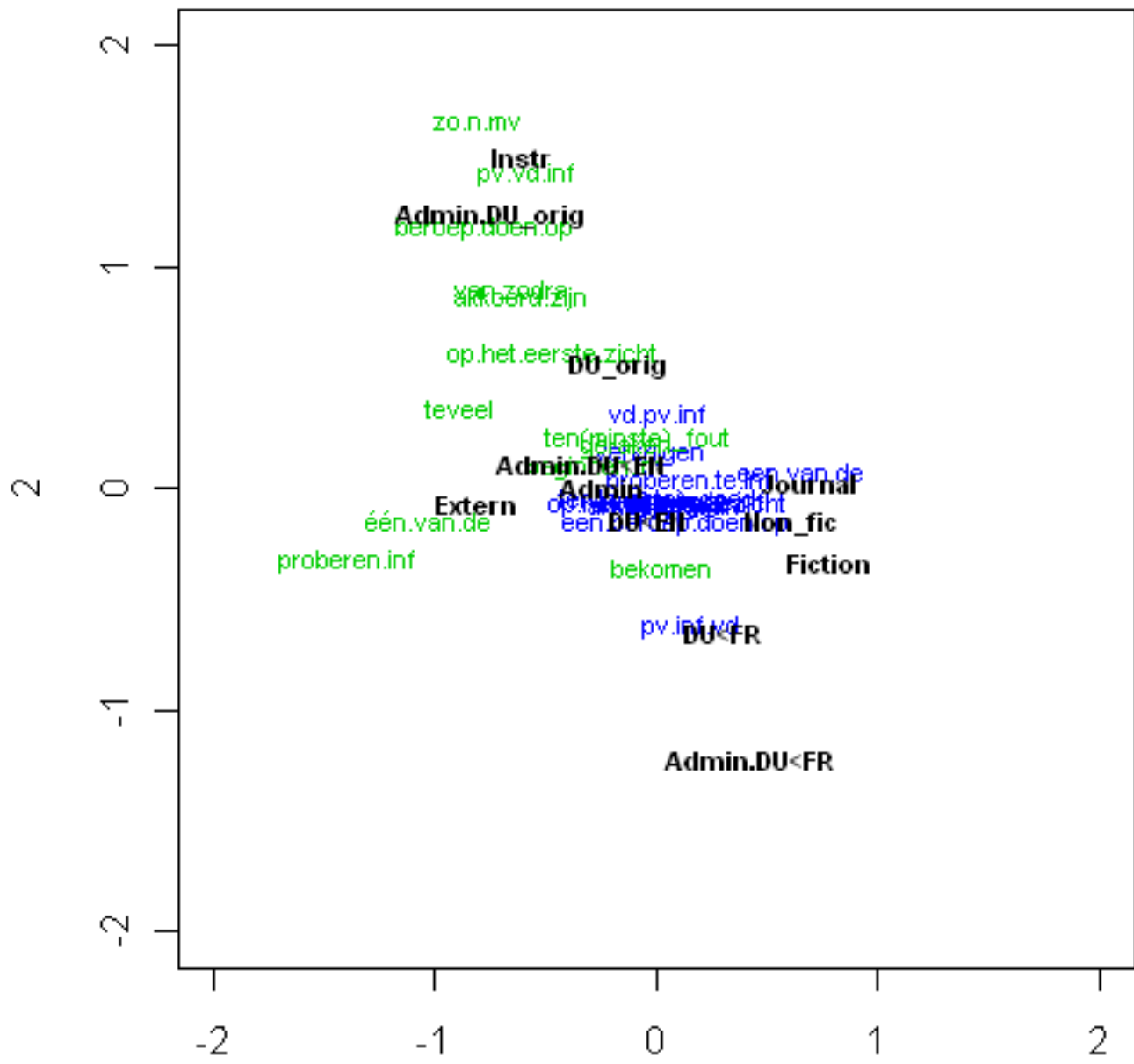
	Adm	Ext	Fic	Ins	Jour	Non_f	DU_or	DU<EN	DU<FR	...
verkrijgen	187	144	1	33	28	17	195	89	126	
bekomen	91	69	0	23	3	6	70	31	91	
zodra	47	48	10	17	85	26	93	37	103	
van.zodra	7	4	0	3	5	0	13	1	5	
pv.inf.vd	63	20	1	6	10	11	33	19	59	
vd.pv.inf	20	14	1	9	19	1	29	16	19	
pv.vd.inf	14	7	0	9	3	0	27	3	3	
...										

'Multifactorial' CA

...	Adm_or	Adm<EN	Adm<FR	Ext_or	Ext<EN	Ext<FR	Fic_or	Fic<EN	Fic<FR	...
	91	48	48	58	38	48	0	0	1	
	25	11	55	19	20	30	0	0	0	
	17	3	27	17	11	20	0	0	10	
	7	0	0	0	1	3	0	0	0	
	14	9	40	2	5	13	0	0	1	
	8	6	6	5	6	3	0	0	1	
	12	1	1	5	2	0	0	0	0	

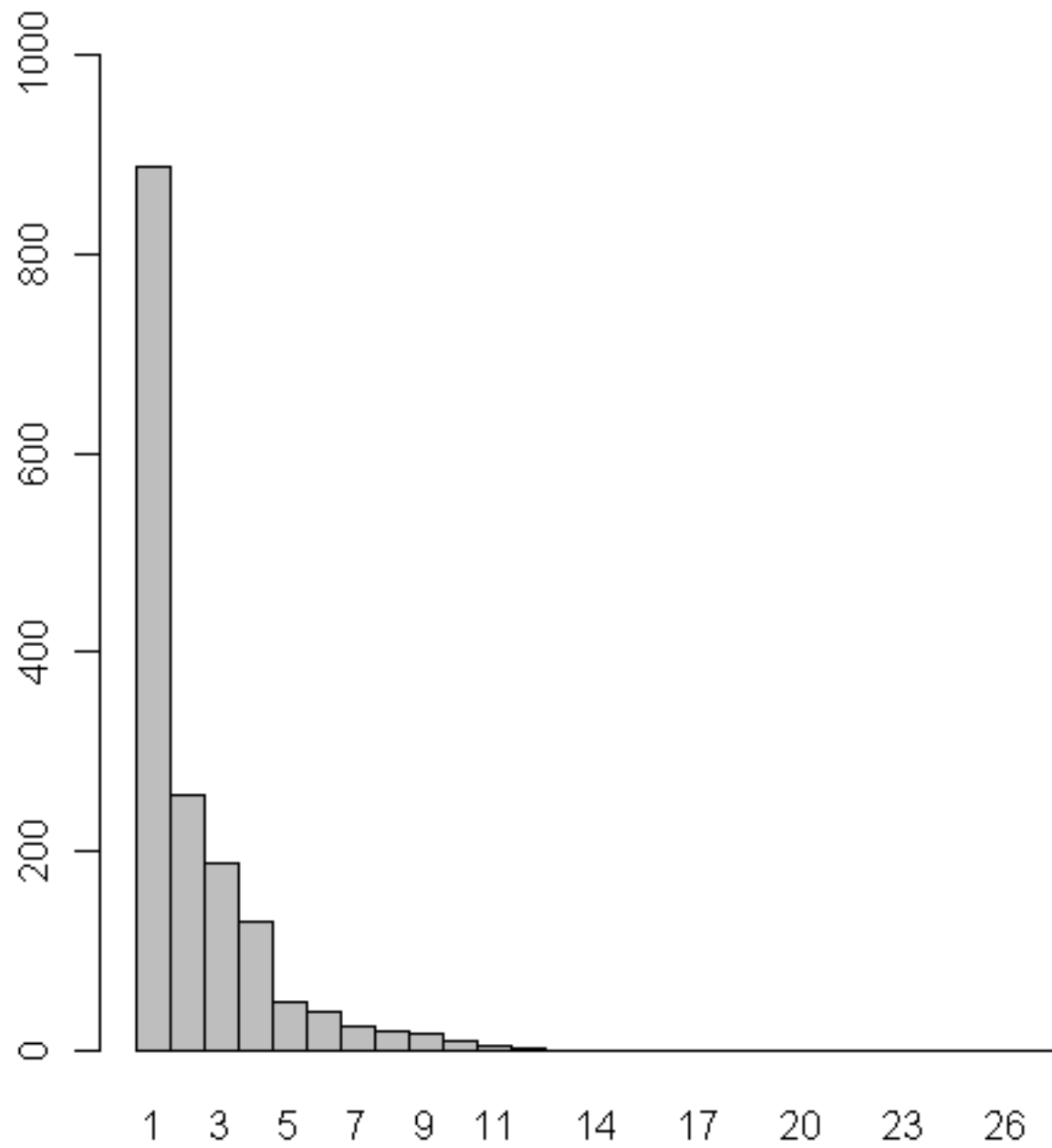
'Multifactorial' CA

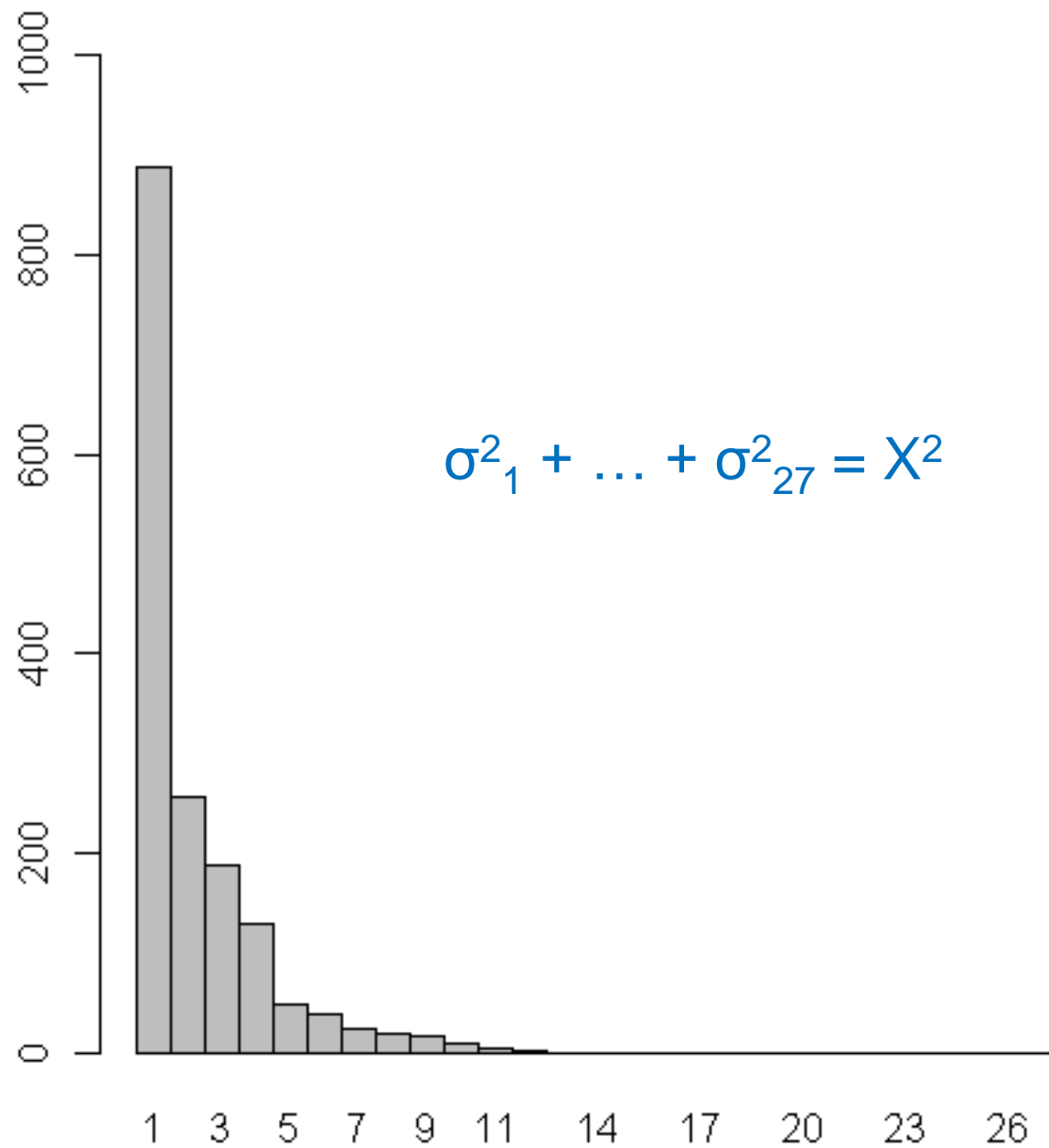
...	Ins_or	Ins<EN	Ins<FR	Jour_or	Jour<EN	Jour<FR	Non_or	Non<EN	Non<FR
	19	0	14	11	3	14	16	0	1
	18	0	5	3	0	0	5	0	1
	11	0	6	28	23	34	20	0	6
	3	0	0	3	0	2	0	0	0
	5	0	1	2	5	3	10	0	1
	7	0	2	8	4	7	1	0	0
	7	0	2	3	0	0	0	0	0

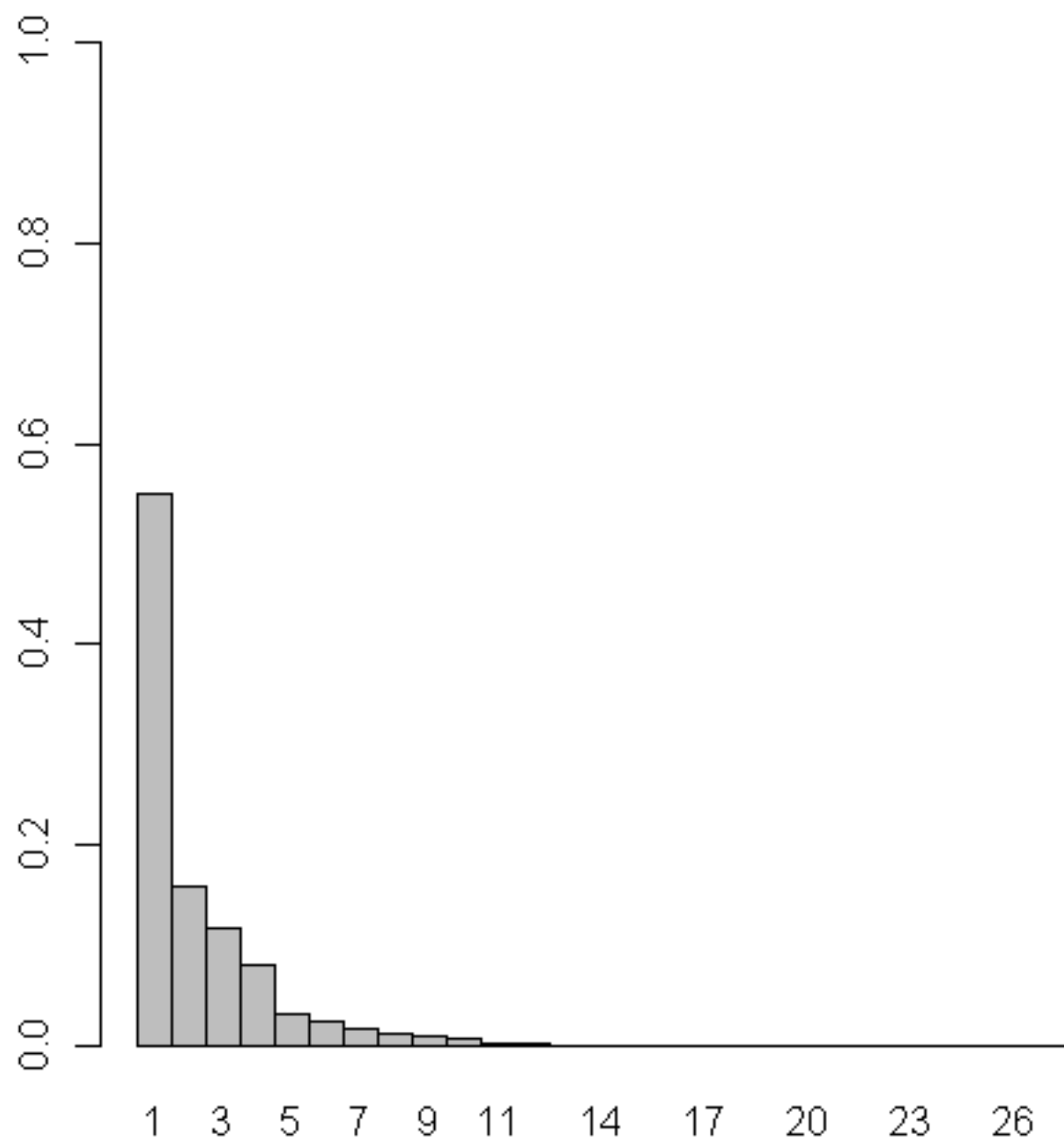


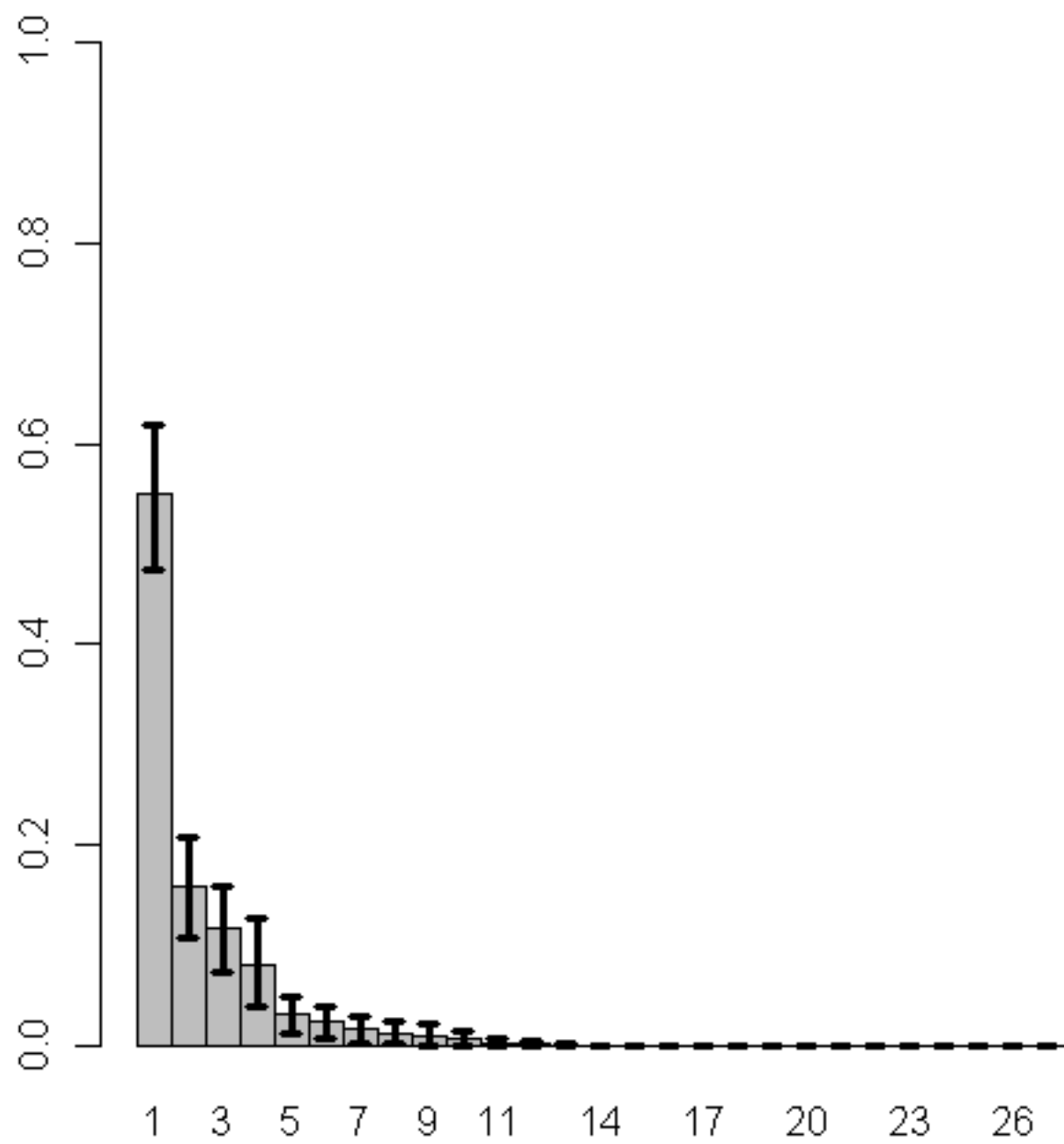
Statistical inference

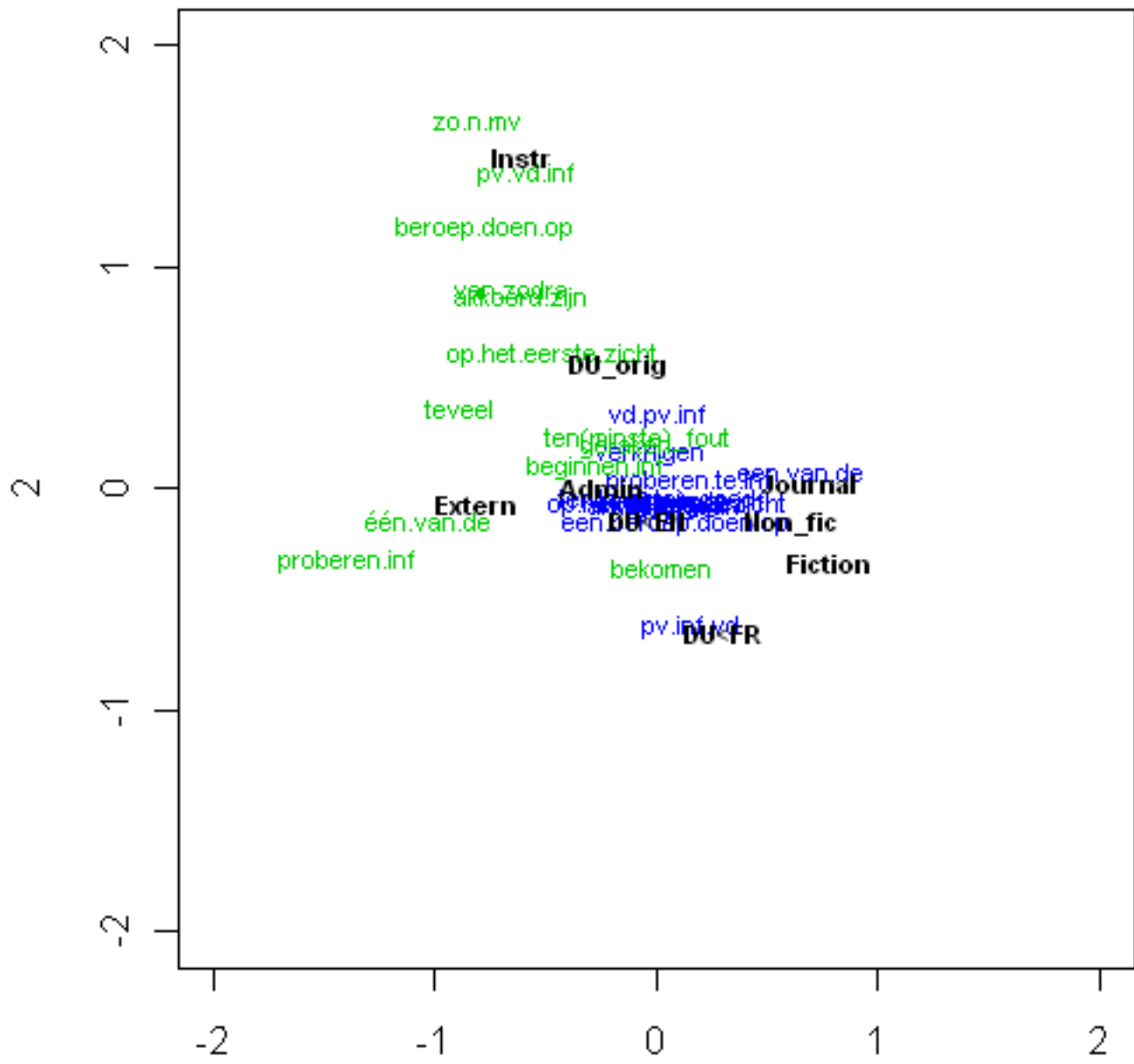
- How many dimensions?
→ Bootstrap confidence intervals for eigenvalues
- Distances between varieties significant?
→ Bootstrap confidence regions for coordinates

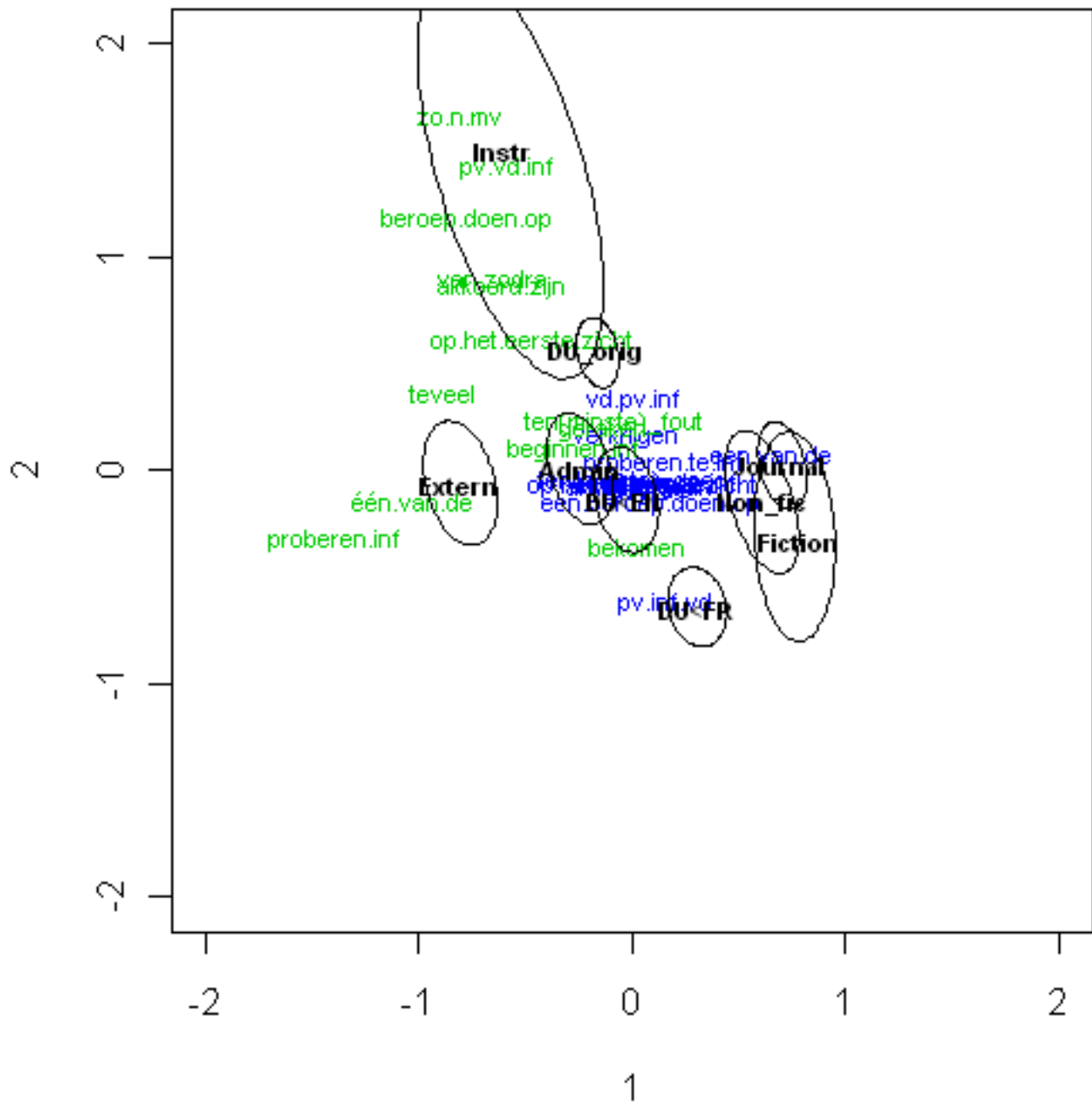


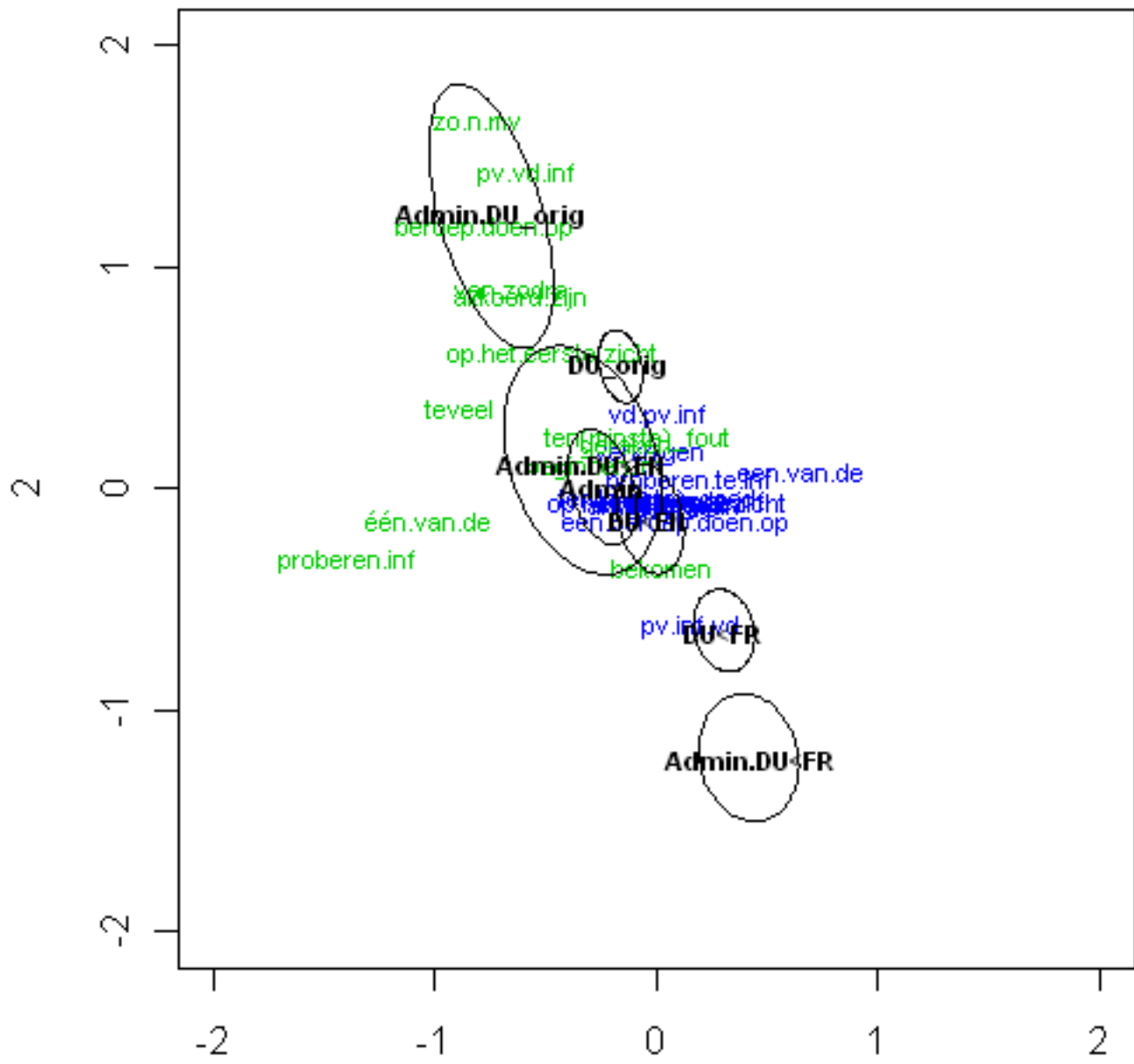


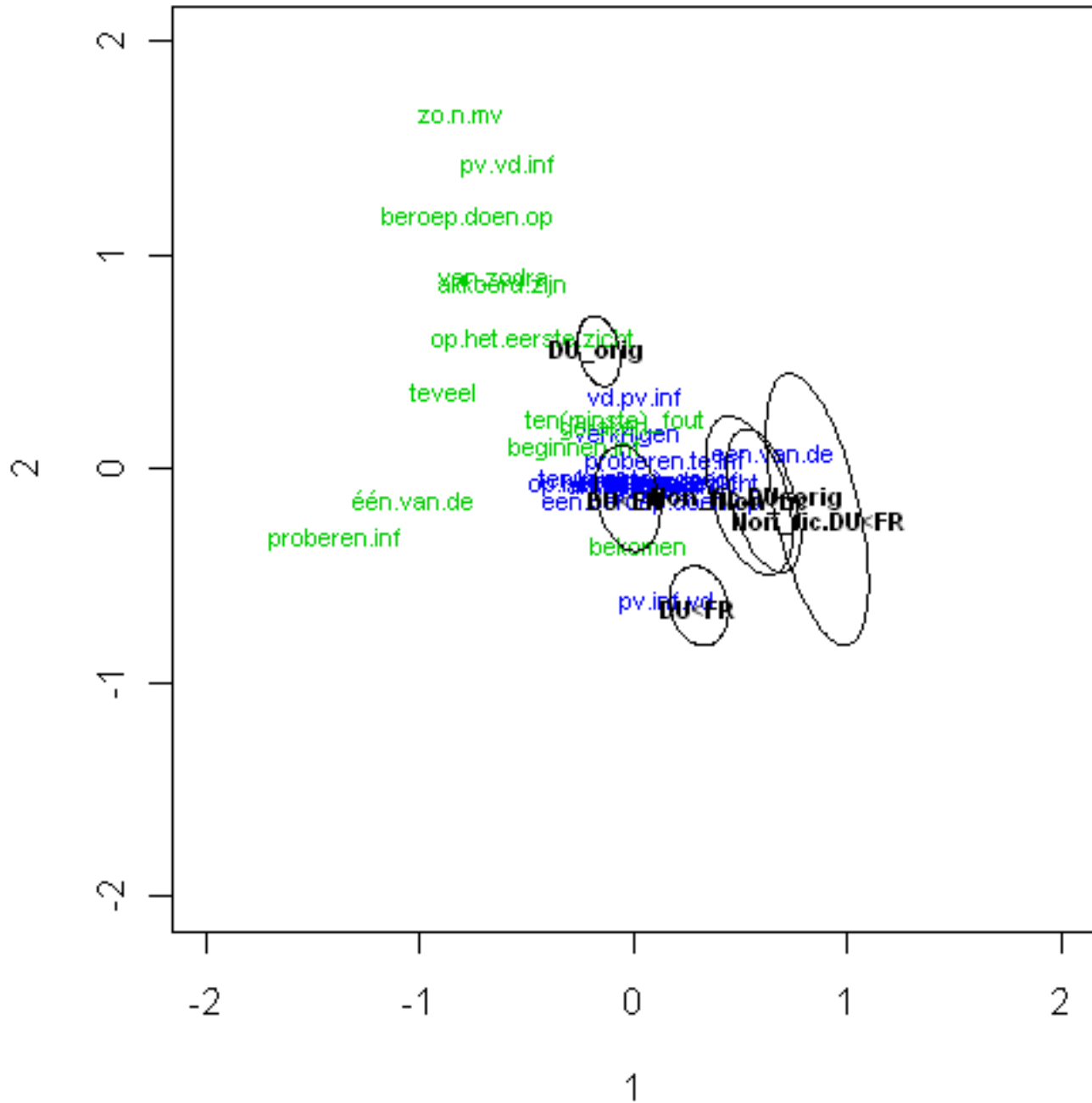










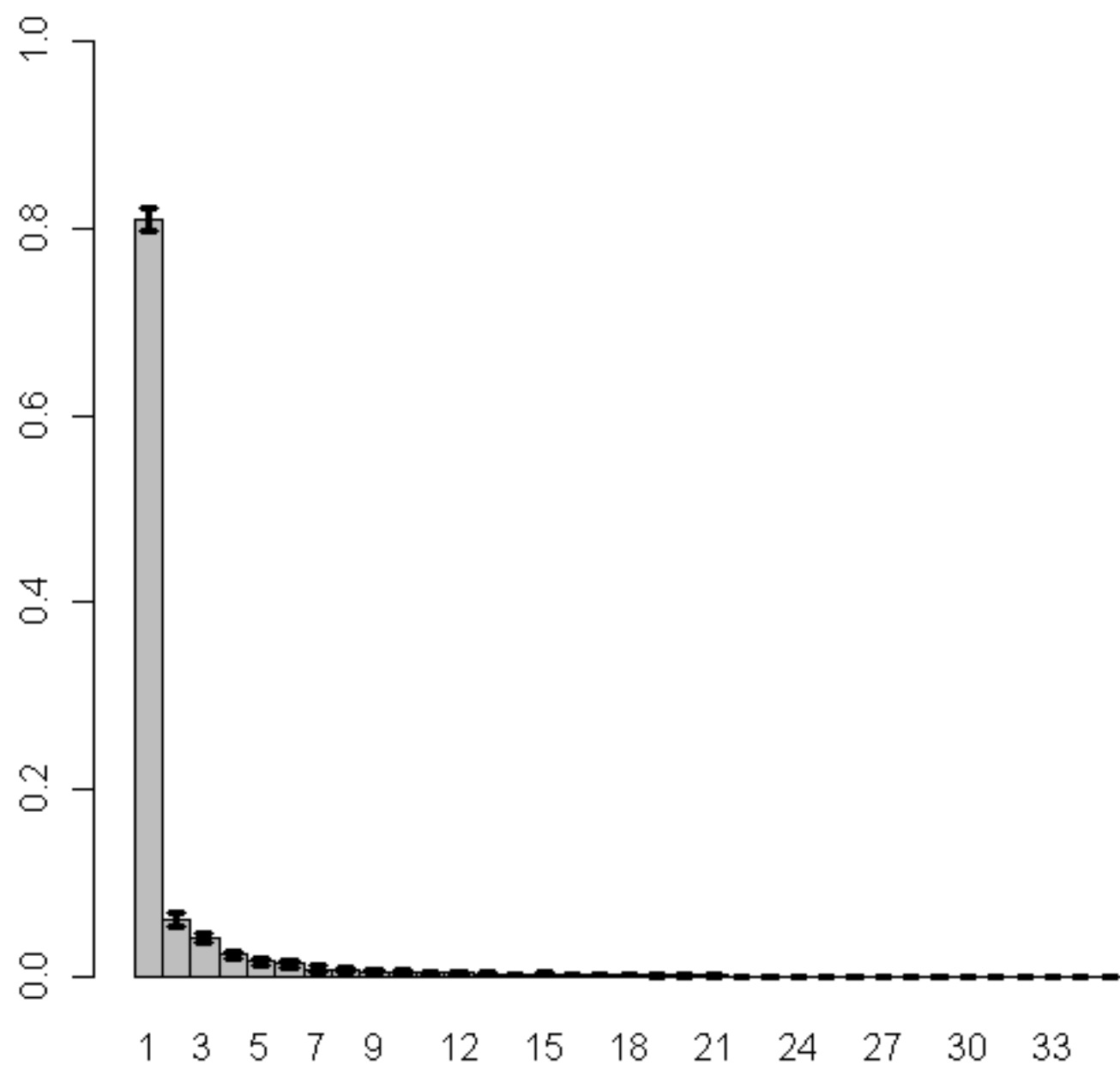


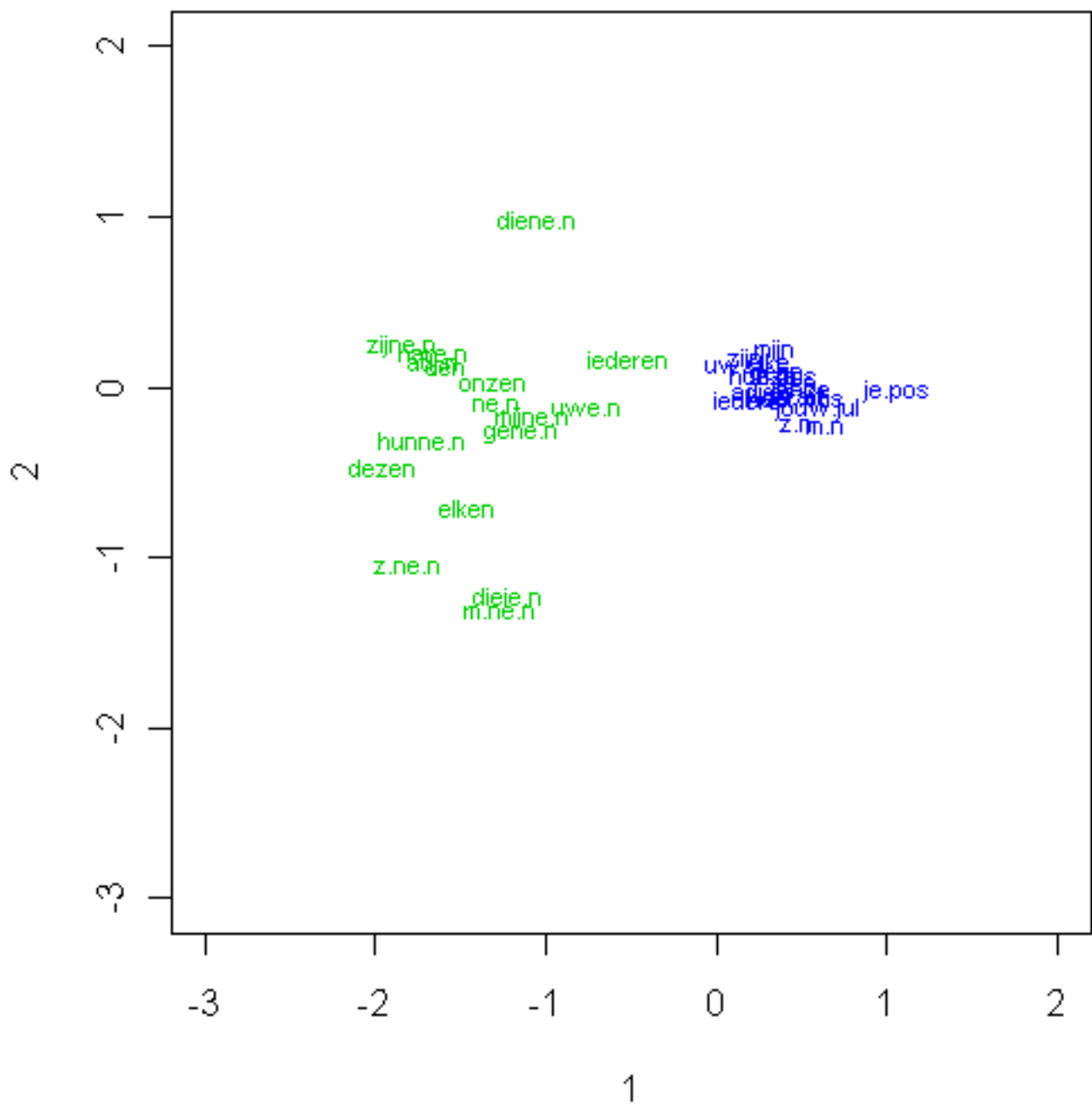
TSS

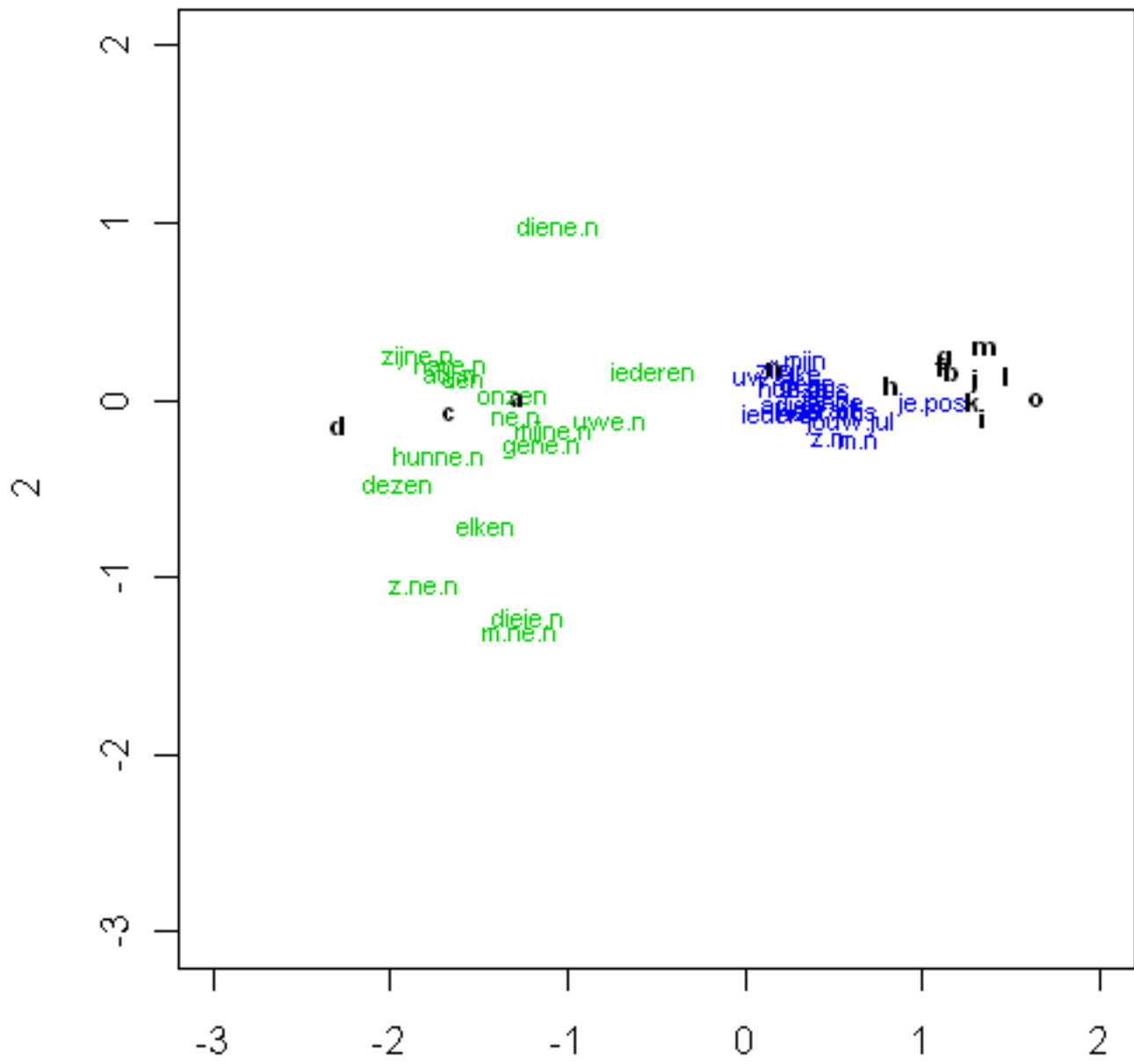
#	Variant 1	Variant 2	Variant 3	Variant 4	Translation
1	de	den			<i>the</i>
2	een	ne(n)			<i>a</i>
3	geen	gene(n)			<i>no</i>
4	elke	elken			<i>each</i>
5	iedere	iederem			<i>every</i>
6	deze	dezen			<i>this</i>
7	die	dieje(n)	diene(n)		<i>that</i>

TSS

#	Variant 1	Variant 2	Variant 3	Variant 4	Translation
8	mijn	m'n	mijne(n)	m'ne(n)	<i>my</i>
9	je	jouw/jullie	uw	uwe(n)	<i>your</i>
10	zijn	z'n	zijne(n)	z'ne(n)	<i>his</i>
11	haar	hare(n)			<i>her</i>
12	onze	onzen			<i>our</i>
13	hun	hunne(n)			<i>their</i>
14	ADJ-e	ADJ-en			<i>adjectival inflection</i>





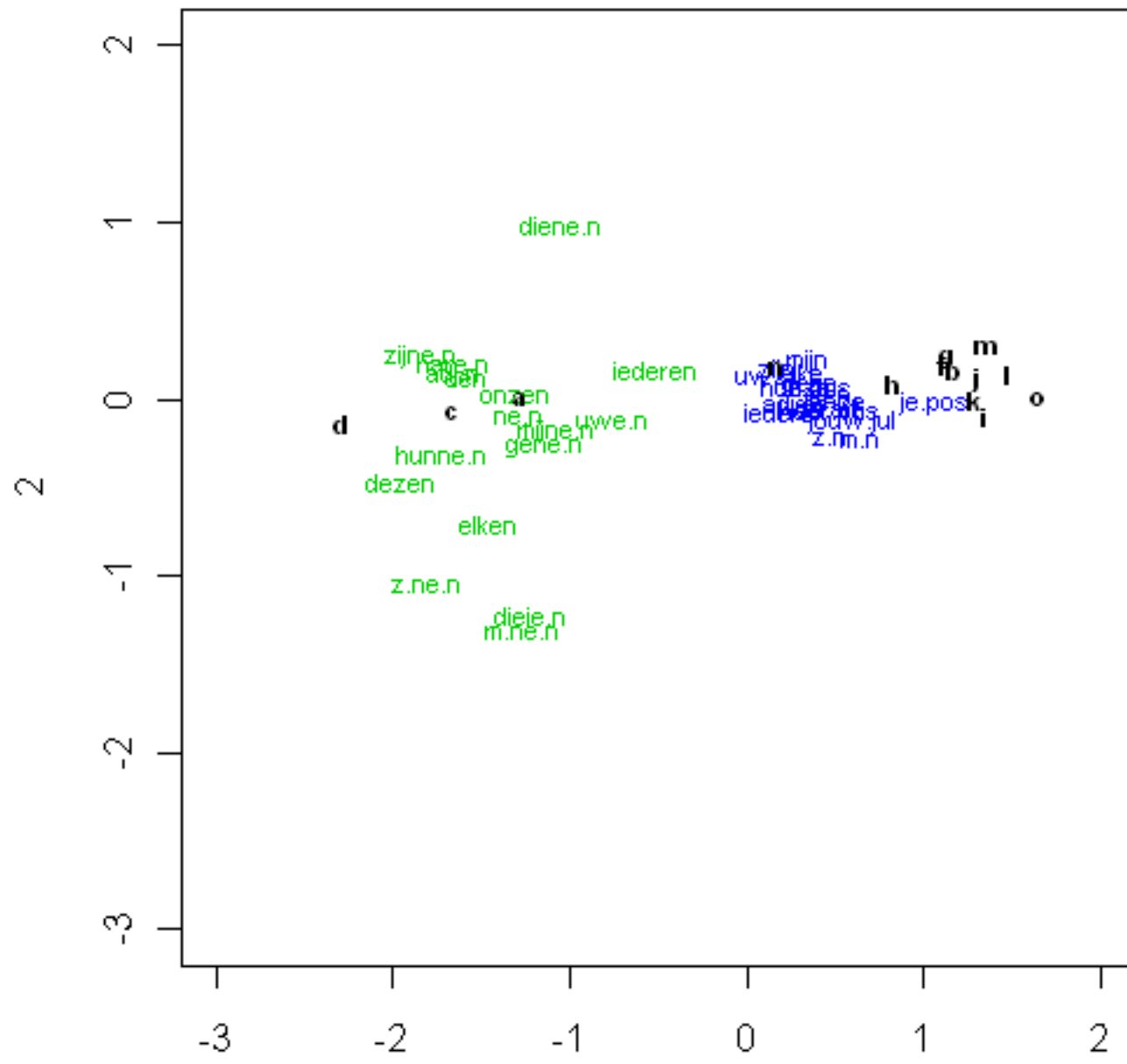


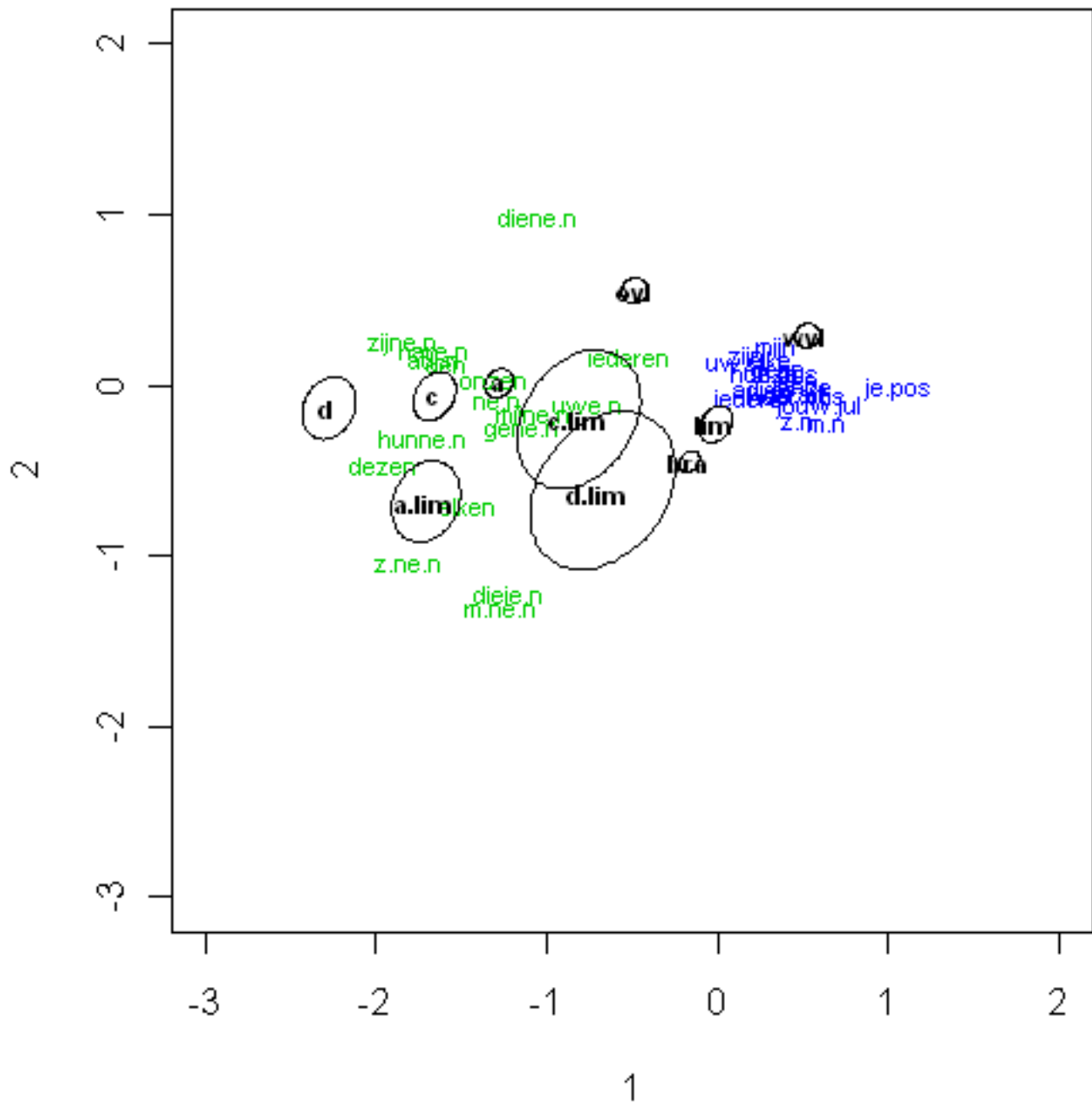
'Components'

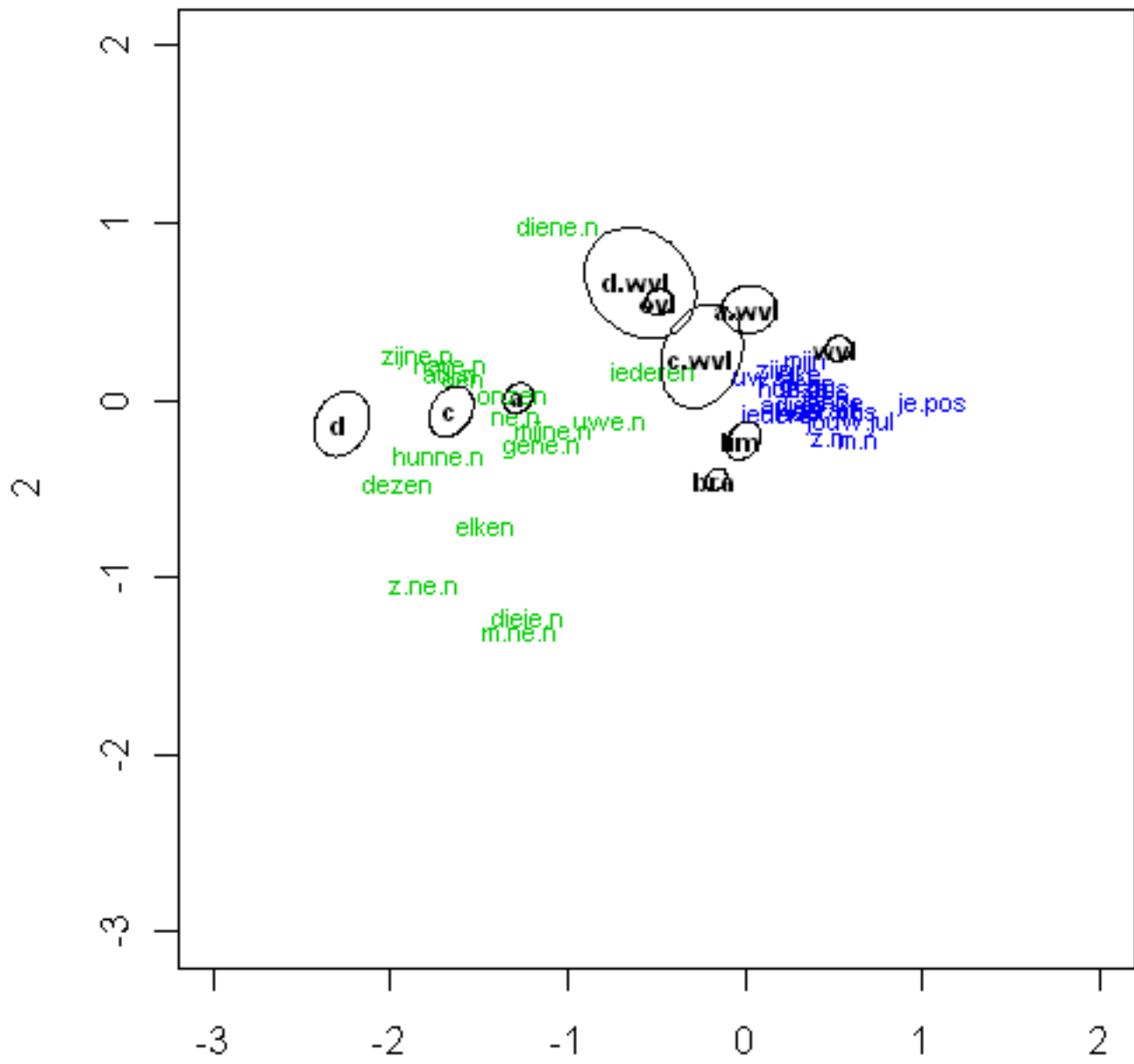
- **a:** Face-to-face conversations
- **b:** Interviews with teachers of Dutch
- **c:** Telephone dialogues (switchboard)
- **d:** Telephone dialogues (mini disc)
- **f:** Broadcast discussions/debates
- **g:** Non-broadcast discussions/debates
- **h:** Lessons recorded in classroom

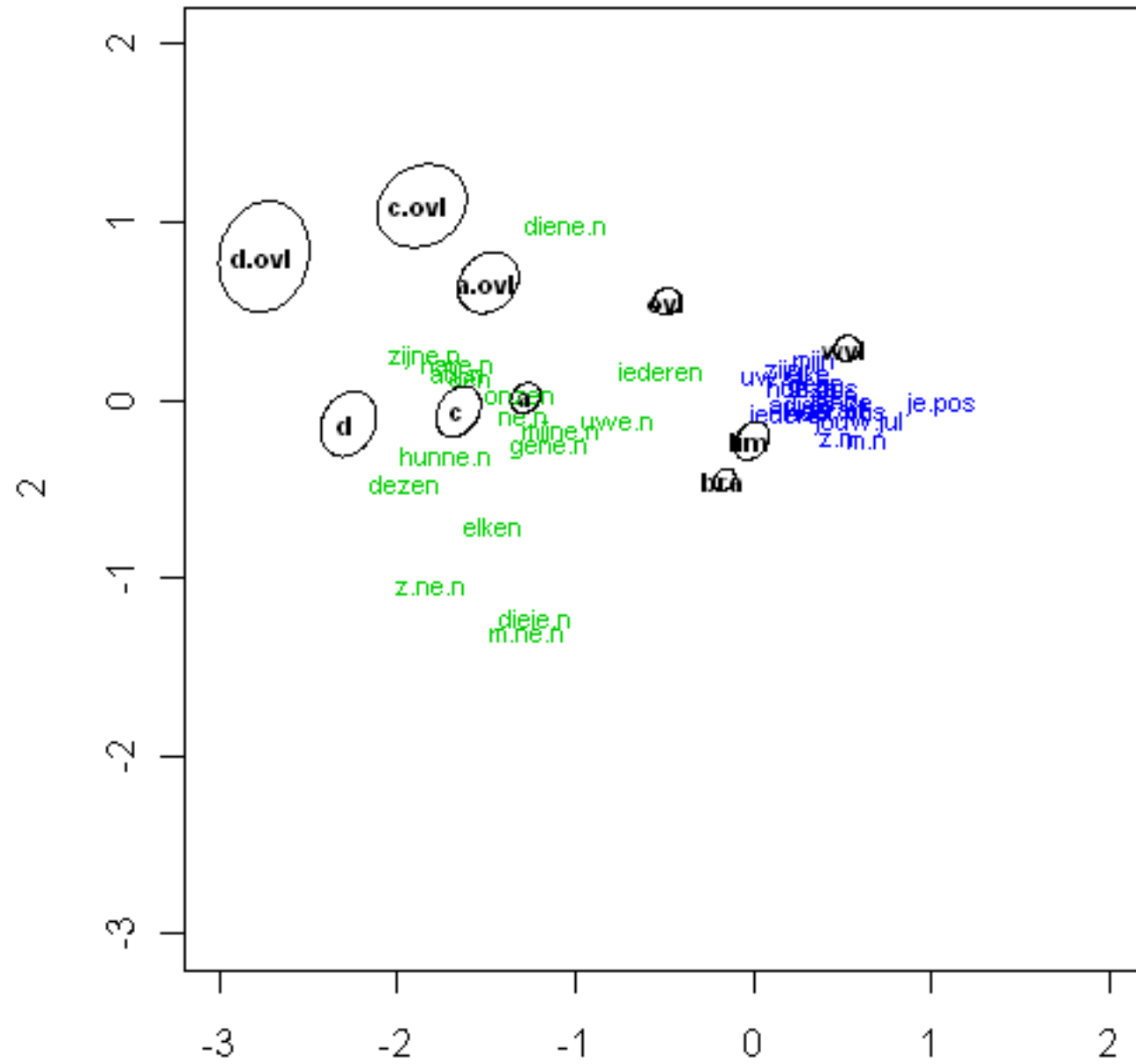
'Components'

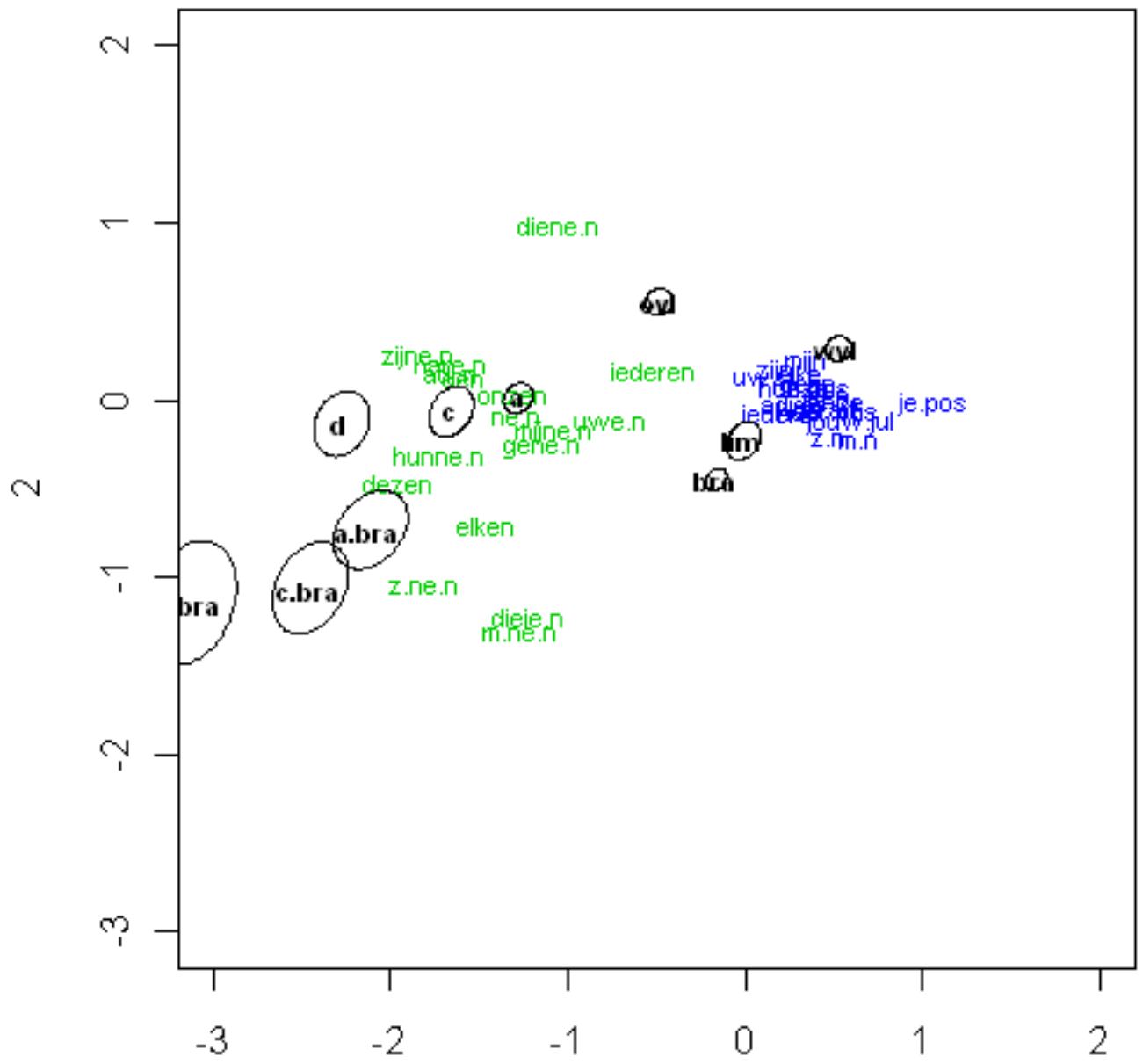
- **i**: Live commentaries (sports)
- **j**: News reports/reportages
- **k**: News bulletins
- **l**: Commentaries/columns/reviews
- **m**: Ceremonious speeches/sermons
- **n**: Lectures/seminars
- **o**: Read texts











Conclusions

- (In)formality depends on **context**: different situations ask for different linguistic items
- (In)formality depends on **power**: dominant individuals can afford themselves more leniency

Thank you!

koen.plevoets@hogent.be