



FRIEDRICH-ALEXANDER  
UNIVERSITÄT  
ERLANGEN-NÜRNBERG

PHILOSOPHISCHE FAKULTÄT  
UND FACHBEREICH THEOLOGIE

# Unit 7: A multivariate approach to linguistic variation

Statistics for Linguists with R – A SIGIL Course

Stephanie Evert

Computational Corpus Linguistics Group  
FAU Erlangen-Nürnberg

# Linguistic variation

## Variation of a quantitative linguistic feature

- frequency of passive, past perfect, split infinitive, ...
- frequency of expression, semantic field, topic, ...
- association strength, lexical density, productivity, ...

## across

- languages and language varieties
- regions & social strata
- time (diachronic change)
- individual speakers & discourses

# Studying linguistic variation

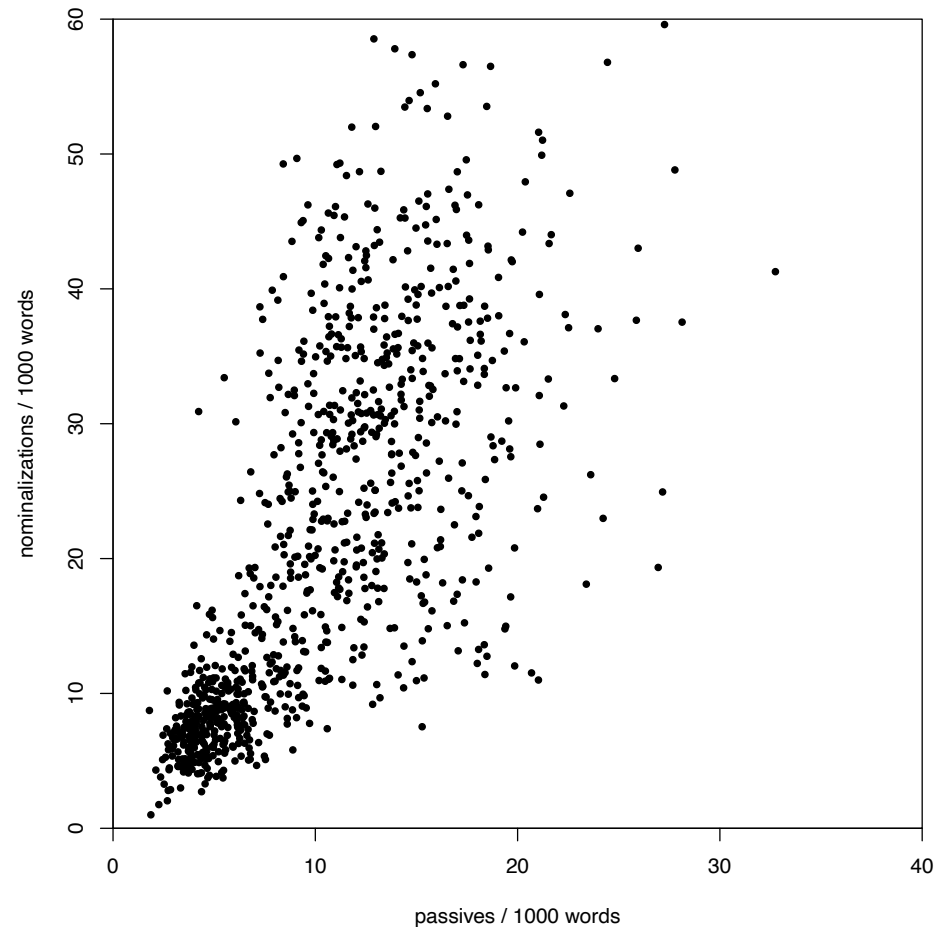
- Univariate approach
  - compare single feature across two or more conditions
  - e.g. AmE vs. BrE vs. IndE vs. ... / male vs. female / etc.
  - corpus frequency comparison
- Regression approach
  - predict single quantity from multiple explanatory factors
- Multivariate approach
  - identify common patterns of variation across multiple different features → correlation analysis
  - inductive techniques don't require pre-defined conditions

# Variation as a nuisance parameter

- Many aspects of linguistic variation are **nuisance parameters** in corpus linguistics
  - e.g. difference in frequency of passives between AmE and BrE, as well as development from 1960s to 1990s (Unit #2)
  - ignore other dimensions such as genre/register variation by **pooling** frequency data from all texts of each corpus
  - corpus is analyzed as a **random sample** of VP tokens
- Consequences
  - variation → non-randomness → overestimate significance
  - discussed in much more detail in Unit #8

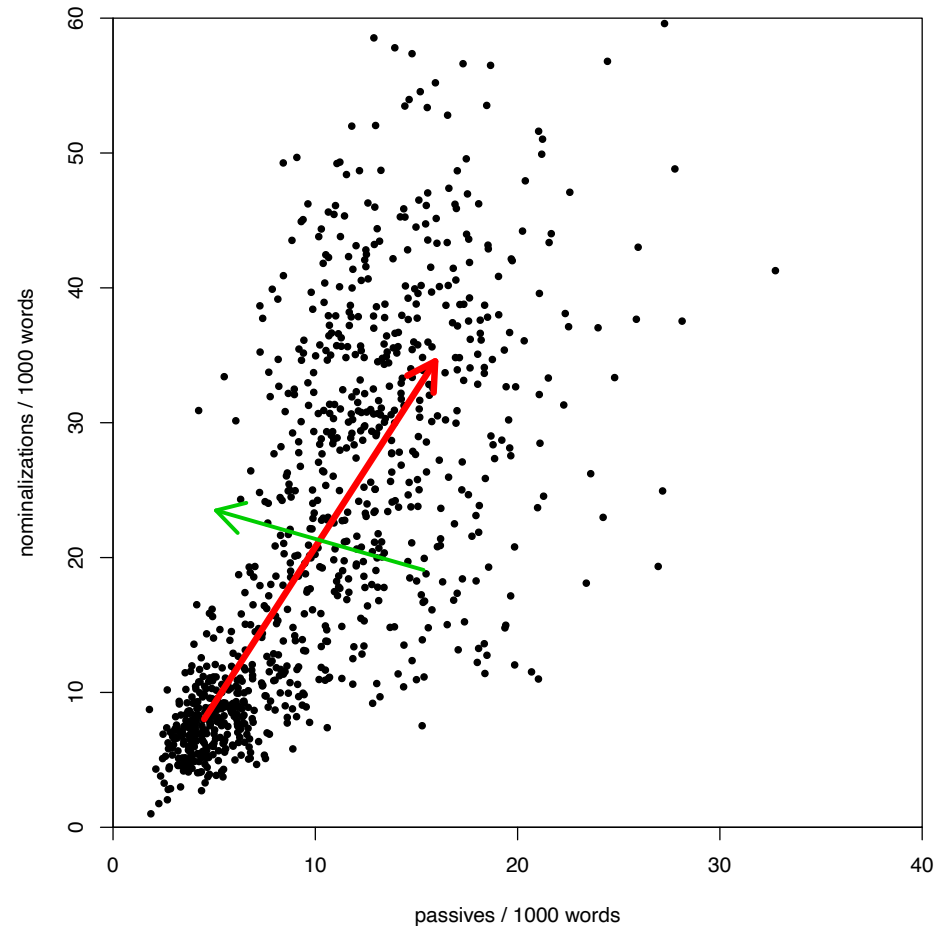
# The multivariate approach

- Different linguistic features often show similar patterns of variation
- E.g. passives and nominalizations

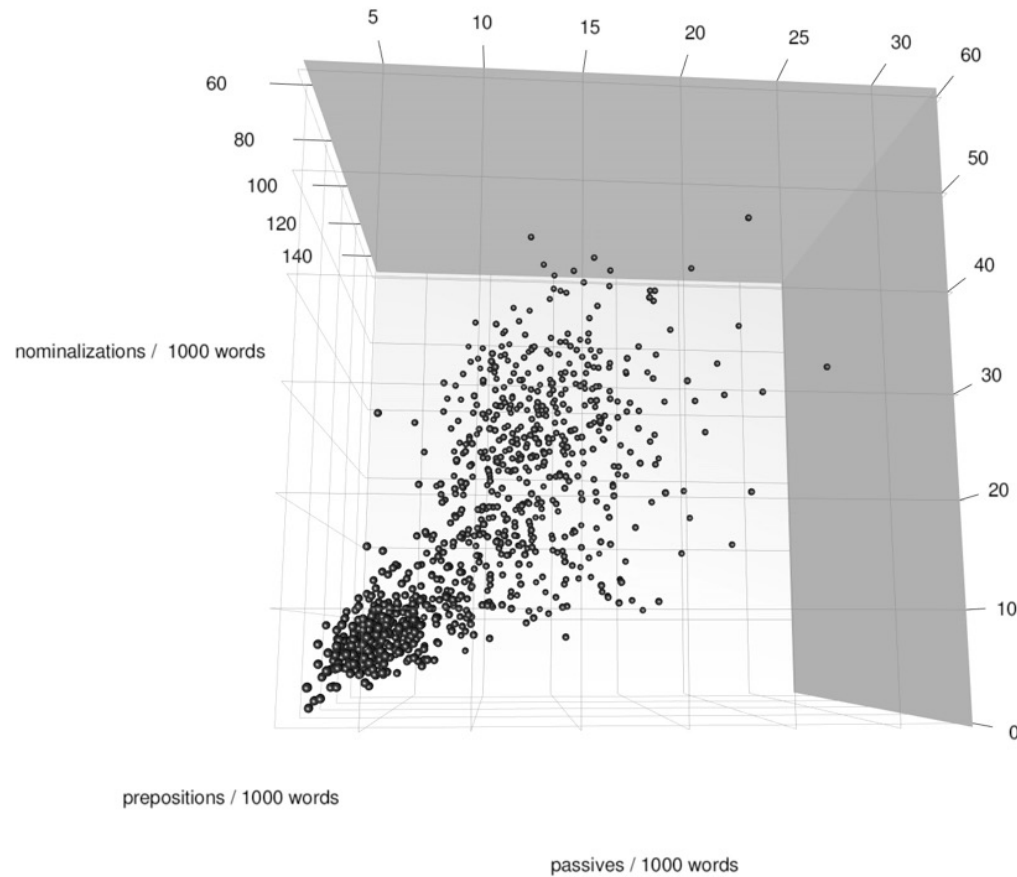


# The multivariate approach

- Different linguistic features often show similar patterns of variation
- E.g. passives and nominalizations
- Such **correlations** can be exploited to determine major **dimensions** of var.



# The multivariate approach



# The multivariate approach

- Multivariate analysis exploits correlations between features in order to determine **latent dimensions**
  - interpreted as underlying “causes” of variation
- An inductive, data-driven approach
  - no theoretical assumptions about linguistic variation and categories / sub-corpora to be compared
- Pioneering work by Doug Biber (1988, 1993, 1995, ...)
  - “multidimensional analysis” of register variation
- Related approaches: correspondence analysis, distributional semantics, topic modelling, ...



# Biber's multidimensional analysis (MDA)

Table 5.7 *Linguistic features used in the analysis of English*


---



---

A. Tense and aspect markers
1 Past tense
2 Perfect aspect
3 Present tense
B. Place and time adverbials
4 Place adverbials (e.g., <i>above, beside, outdoors</i> )
5 Time adverbials (e.g., <i>early, instantly, soon</i> )
C. Pronouns and pro-verbs
6 First-person pronouns
7 Second-person pronouns
8 Third-person personal pronouns (excluding <i>it</i> )
9 Pronoun <i>it</i>
10 Demonstrative pronouns ( <i>that, this, these, those</i> as pronouns)
11 Indefinite pronouns (e.g., <i>anybody, nothing, someone</i> )
12 Pro-verb <i>do</i>
D. Questions
13 Direct WH questions
E. Nominal forms
14 Nominalizations (ending in <i>-tion, -ment, -ness, -ity</i> )
15 Gerunds (participial forms functioning as nouns)
16 Total other nouns
F. Passives
17 Agentless passives
18 <i>by</i> -passives
G. Stative forms
19 <i>be</i> as main verb
20 Existential <i>there</i>
H. Subordination features
21 <i>that</i> verb complements (e.g., <i>I said that he went.</i> )
22 <i>that</i> adjective complements (e.g., <i>I'm glad that you like it.</i> )
23 WH-clauses (e.g., <i>I believed what he told me.</i> )
24 Infinitives
25 Present participial adverbial clauses (e.g., <i>Stuffing his mouth with cookies, Joe ran out the door.</i> )
26 Past participial adverbial clauses (e.g., <i>Built in a single week, the house would stand for fifty years.</i> )
27 Past participial postnominal (reduced relative) clauses (e.g., <i>the solution produced by this process</i> )
28 Present participial postnominal (reduced relative) clauses (e.g., <i>The event causing this decline was . . .</i> )
29 <i>that</i> relative clauses on subject position (e.g., <i>the dog that bit me</i> )
30 <i>that</i> relative clauses on object position (e.g., <i>the dog that I saw</i> )
31 WH relatives on subject position (e.g., <i>the man who likes popcorn</i> )
32 WH relatives on object position (e.g., <i>the man who Sally likes</i> )
33 Pied-piping relative clauses (e.g., <i>the manner in which he was told</i> )

---



---

Table 5.7 (cont.)

---



---

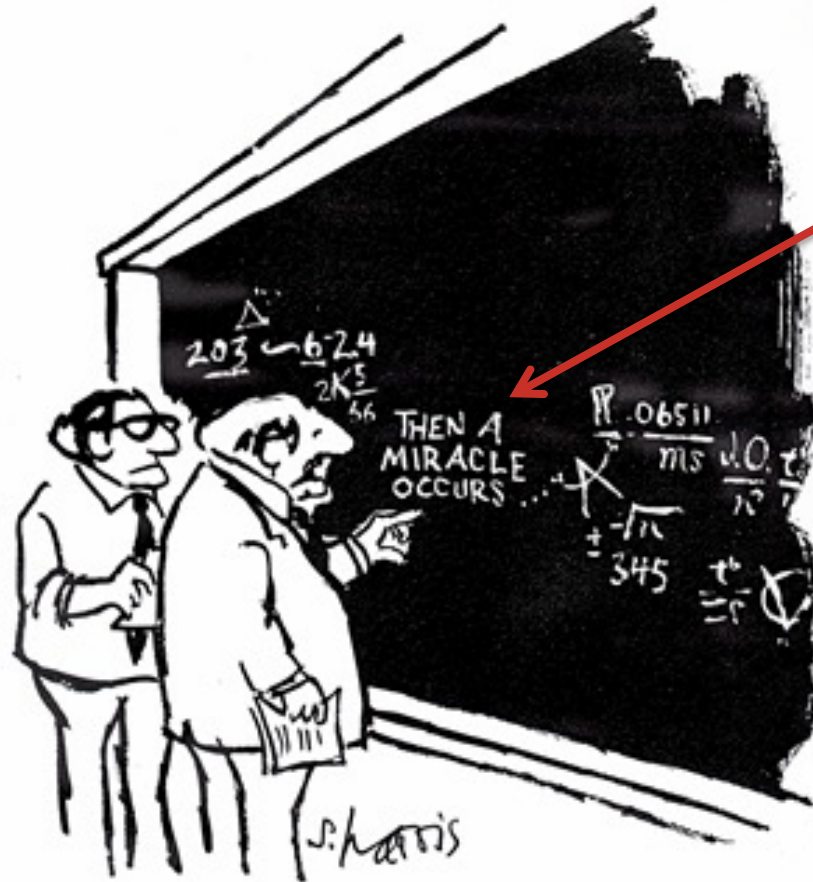
34 Sentence relatives (e.g., <i>Bob likes fried mangoes, which is the most disgusting thing I've ever heard of.</i> )
35 Causative adverbial subordinator ( <i>because</i> )
36 Concessive adverbial subordinators ( <i>although, though</i> )
37 Conditional adverbial subordinators ( <i>if, unless</i> )
38 Other adverbial subordinators (e.g., <i>since, while, whereas</i> )
I. Prepositional phrases, adjectives, and adverbs
39 Total prepositional phrases
40 Attributive adjectives (e.g., <i>the big horse</i> )
41 Predicative adjectives (e.g., <i>The horse is big.</i> )
42 Total adverbs
J. Lexical specificity
43 Type-token ratio
44 Mean word length
K. Lexical classes
45 Conjunctions (e.g., <i>consequently, furthermore, however</i> )
46 Downtoners (e.g., <i>barely, nearly, slightly</i> )
47 Hedges (e.g., <i>at about, something like, almost</i> )
48 Amplifiers (e.g., <i>absolutely, extremely, perfectly</i> )
49 Emphatics (e.g., <i>a lot, for sure, really</i> )
50 Discourse particles (e.g., sentence-initial <i>well, now, anyway</i> )
51 Demonstratives
L. Modals
52 Possibility modals ( <i>can, may, might, could</i> )
53 Necessity modals ( <i>ought, should, must</i> )
54 Predictive modals ( <i>will, would, shall</i> )
M. Specialized verb classes
55 Public verbs (e.g., <i>assert, declare, mention</i> )
56 Private verbs (e.g., <i>assume, believe, doubt, know</i> )
57 Suasive verbs (e.g., <i>command, insist, propose</i> )
58 <i>seem</i> and <i>appear</i>
N. Reduced forms and dispreferred structures
59 Contractions
60 Subordinator <i>that</i> deletion (e.g., <i>I think [that] he went.</i> )
61 Stranded prepositions (e.g., <i>the candidate that I was thinking of</i> )
62 Split infinitives (e.g., <i>He wants to convincingly prove that . . .</i> )
63 Split auxiliaries (e.g., <i>They were apparently shown to . . .</i> )
O. Co-ordination
64 Phrasal co-ordination (NOUN <i>and</i> NOUN; ADJ; <i>and</i> ADJ; VERB <i>and</i> VERB; ADV <i>and</i> ADV)
65 Independent clause co-ordination (clause-initial <i>and</i> )
P. Negation
66 Synthetic negation (e.g., <i>No answer is good enough for Jones.</i> )
67 Analytic negation (e.g., <i>That's not likely</i> )

---



---

# Biber's MDA



factor analysis  
(FA)

"I THINK YOU SHOULD BE MORE  
EXPLICIT HERE IN STEP TWO."

# Biber's MDA



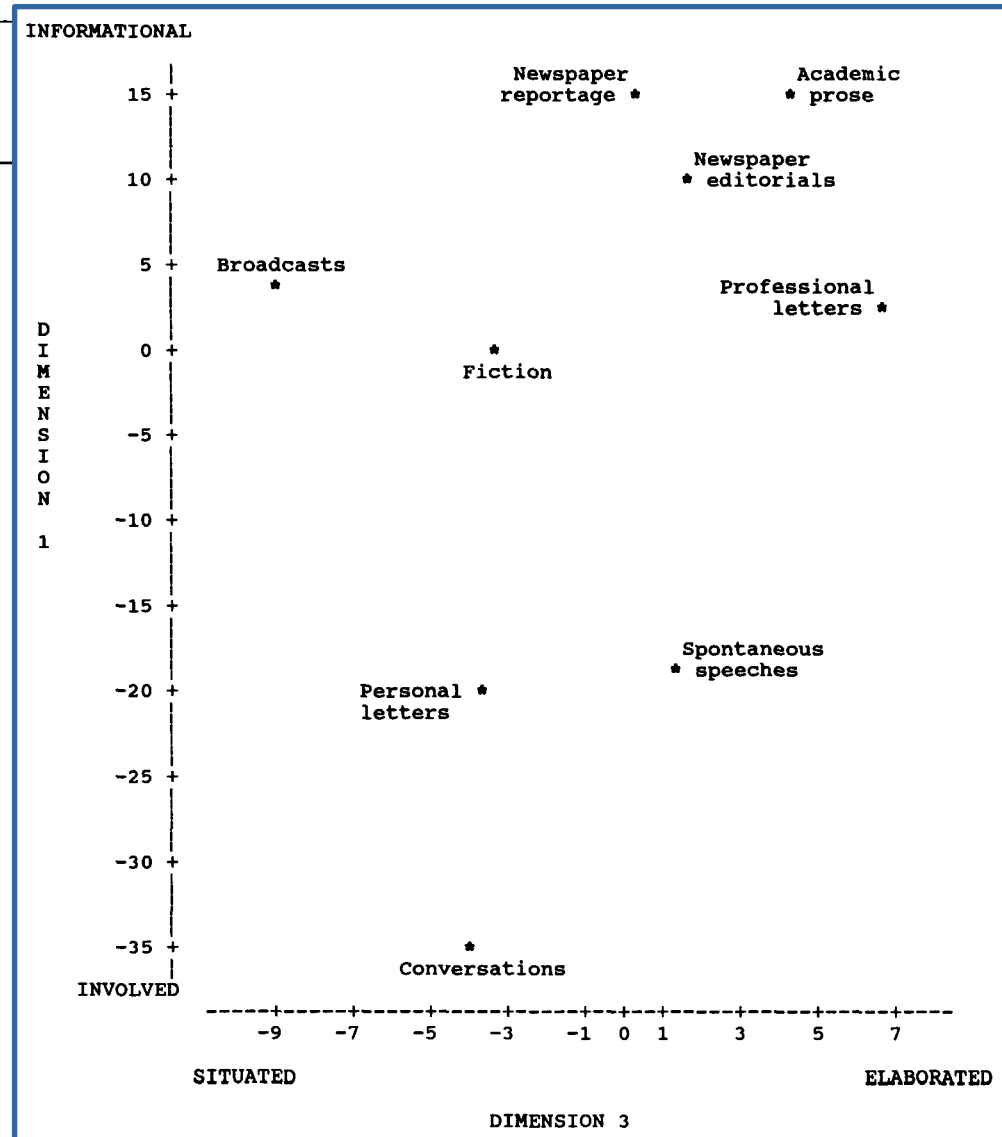
FRIEDRICH-ALEXANDER  
UNIVERSITÄT  
ERLANGEN-NÜRNBERG

PHILOSOPHISCHE FAKULTÄT  
UND FACHBEREICH THEOLOGIE

TABLE 2

Summary of the co-occurrence patterns underlying five major dimensions of English.

DIMENSION 1 (Informational vs. Involved)		DIMENSION 2 (Narrative versus Non-Narrative)	
nouns	0.80	past tense verbs	0.90
word length	0.58	third person pronouns	0.73
prepositional phrases	0.54	perfect aspect verbs	0.48
type / token ratio	0.54	public verbs	0.43
attributive adjs.	0.47	synthetic negation	0.40
private verbs	-0.96	present participial clauses	0.39
<i>that</i> deletions	-0.91	present tense verbs	-0.47
contractions	-0.90	attributive adjs.	-0.41
present tense verbs	-0.86		
2nd person pronouns	-0.86		
<i>do</i> as pro-verb	-0.82		
analytic negation	-0.78		
demonstrative pronouns	-0.76		
general emphatics	-0.74		
first person pronouns	-0.74		
pronoun <i>it</i>	-0.71		
<i>be</i> as main verb	-0.71		
causative subordination	-0.66		
discourse particles	-0.66		
indefinite pronouns	-0.62		
general hedges	-0.58		
amplifiers	-0.56		
sentence relatives	-0.55		
WH questions	-0.52		
possibility modals	-0.50		
non-phrasal coordination	-0.48		
WH clauses	-0.47		
final prepositions	-0.43		



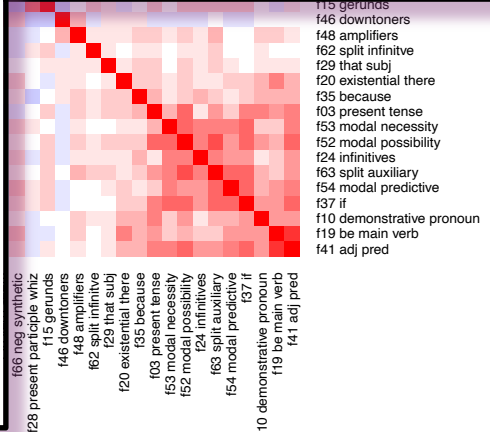
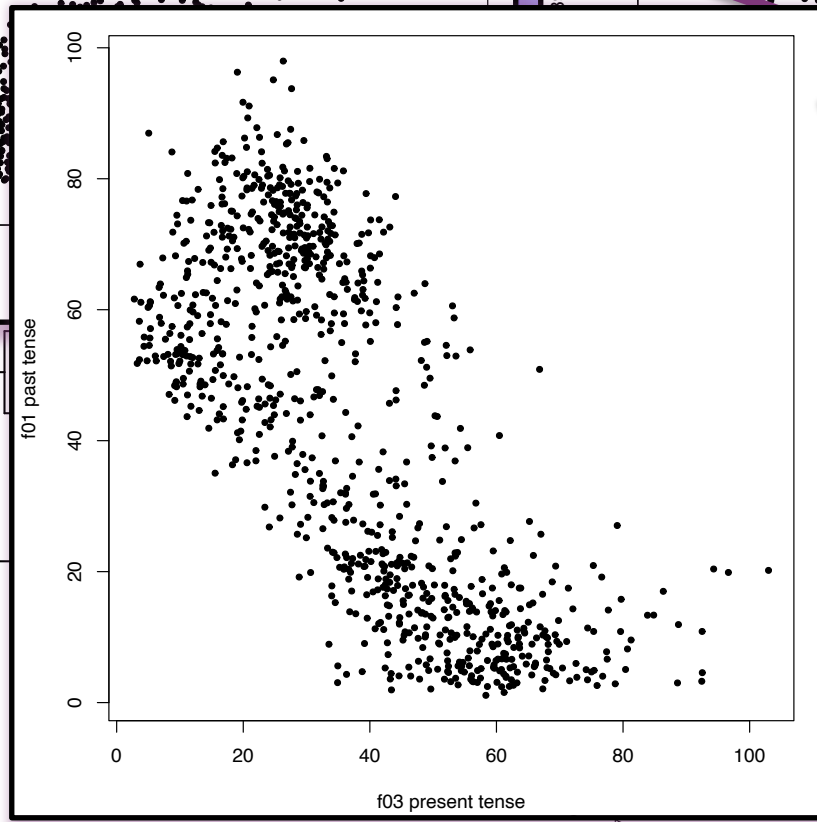
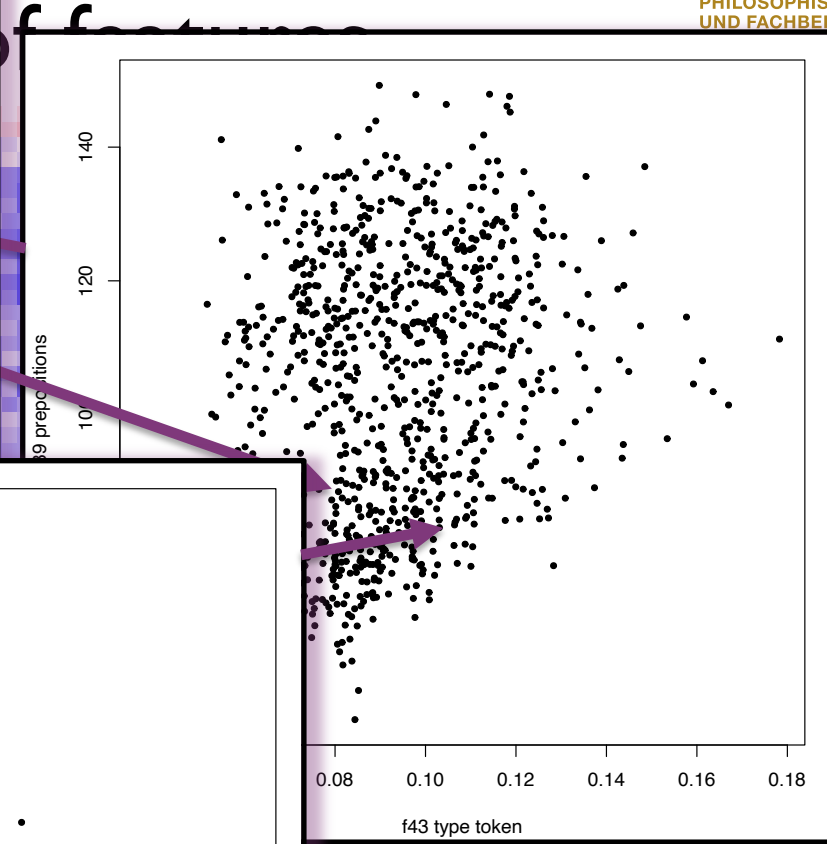
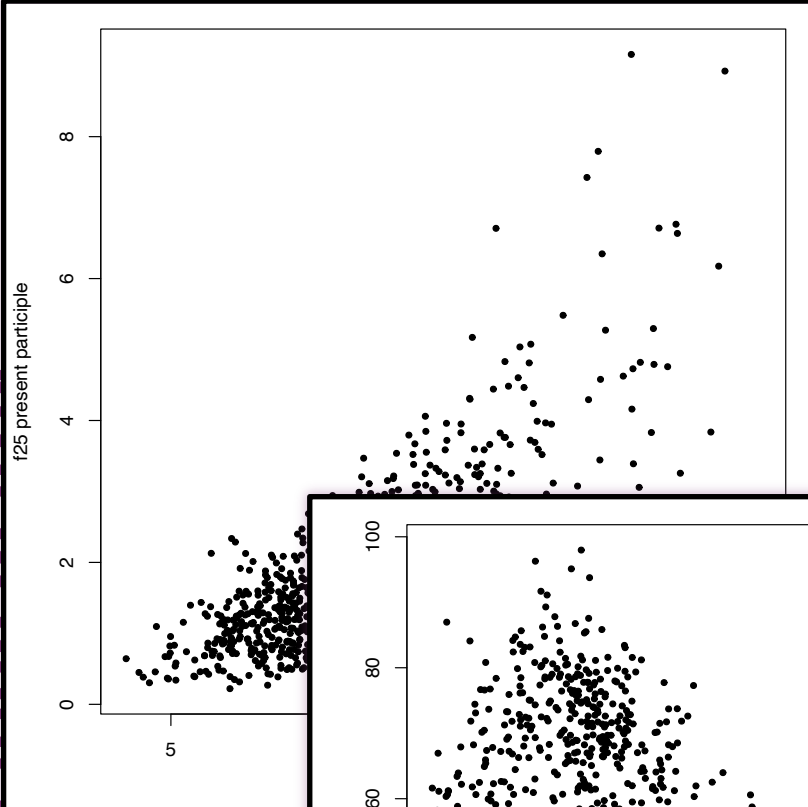
# Pitfalls

- Design bias: choice of quantitative features
- Design bias: selection of text samples
- Involves a miracle
  - not clear what quantitative patterns are captured by FA
  - magic number: how many factor dimensions?
- Interpretation bias
  - arbitrary cutoff for feature weights (“loadings”)
  - risk of reading one's own expectations into features
- More subtle patterns of variation invisible
- Significance & reproducibility of results?

# Reproducing Biber's dimensions

- Sample of 923 medium-length published texts from written part of British National Corpus (BNC)
- Covers 4 different text types + male/female authors
  - academic writing, non-academic prose, fiction, misc.
- Biber features extracted automatically with Python script (Gasthaus 2007)
  - all frequencies normalized per 1000 words
  - data available in R package `corpora` (`BNCbiber`)
- Factor analysis with 4 latent dimensions + varimax
  - seems to yield the most clearly structured analysis

# correlation matrix



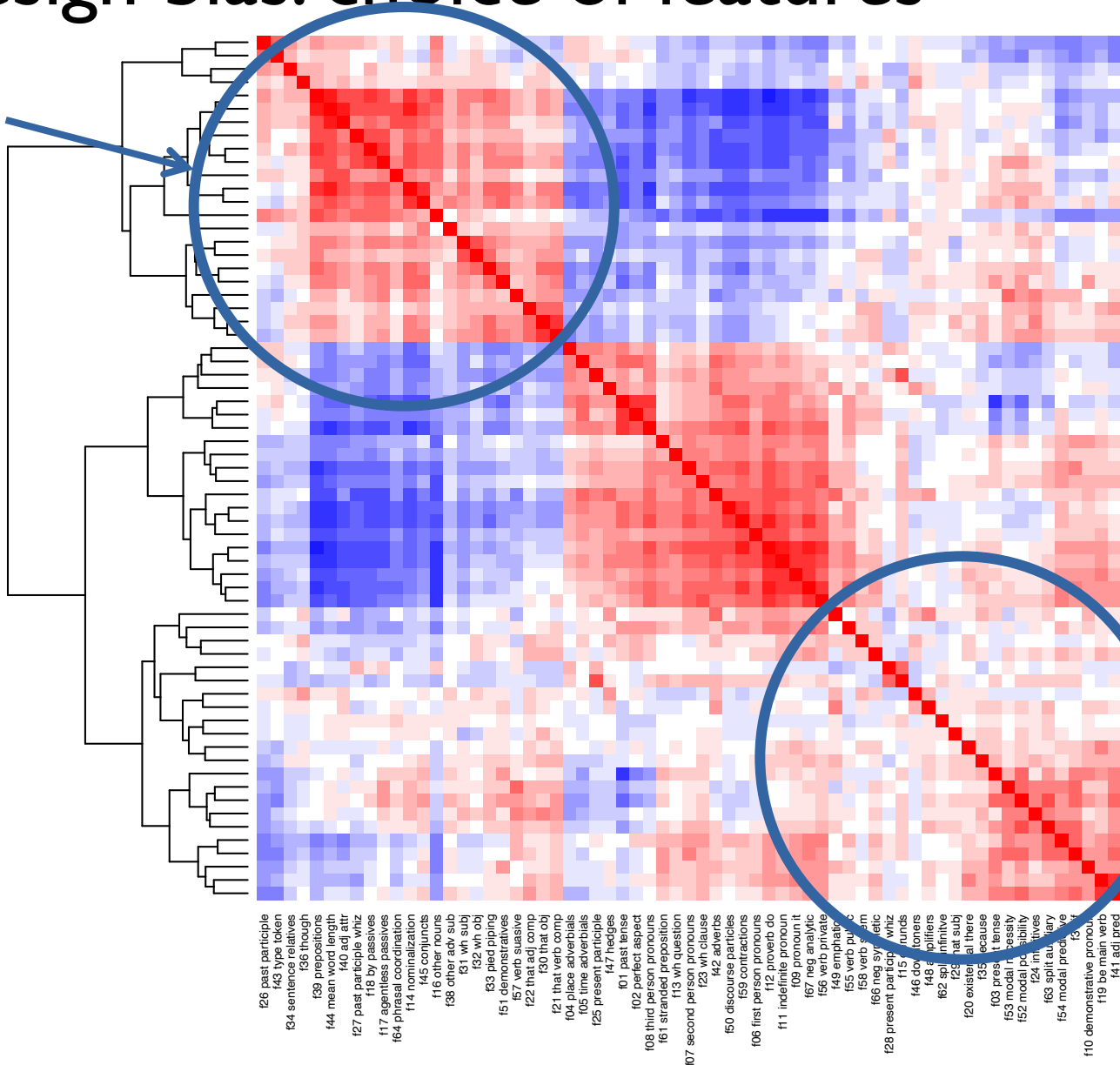
- f66 neg synthetic
- f28 present participle whiz
- f15 gerunds
- f46 downtoners
- f48 amplifiers
- f62 split infinitive
- f29 that subj
- f20 existential there
- f35 because
- f03 present tense
- f53 modal necessity
- f52 modal possibility
- f24 infinitives
- f63 split auxiliary
- f54 modal predictive
- f37 if
- f10 demonstrative pronoun
- f19 be main verb
- f41 adj pred

- f15 gerunds
- f46 downtoners
- f48 amplifiers
- f62 split infinitive
- f29 that subj
- f20 existential there
- f35 because
- f03 present tense
- f53 modal necessity
- f52 modal possibility
- f24 infinitives
- f63 split auxiliary
- f54 modal predictive
- f37 if
- f10 demonstrative pronoun
- f19 be main verb
- f41 adj pred

# Design bias: choice of features

correlation matrix

correlated with noun frequency

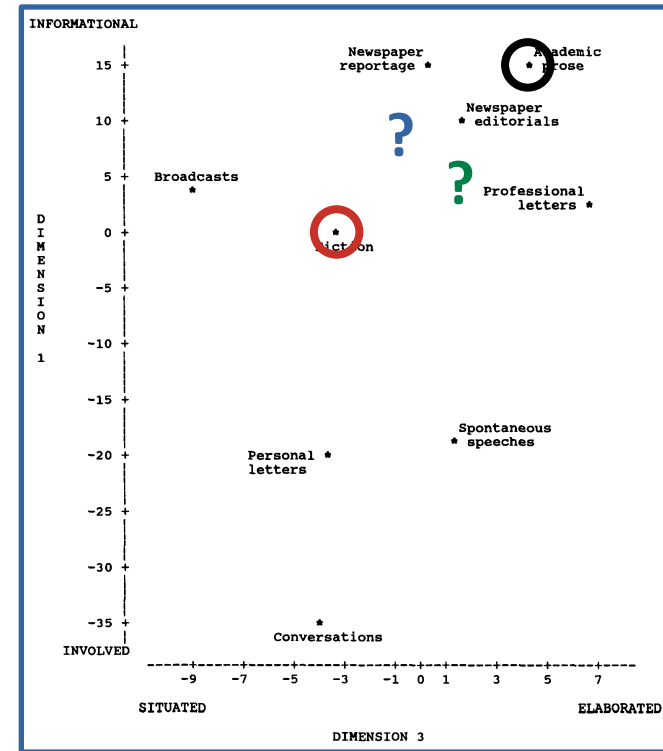
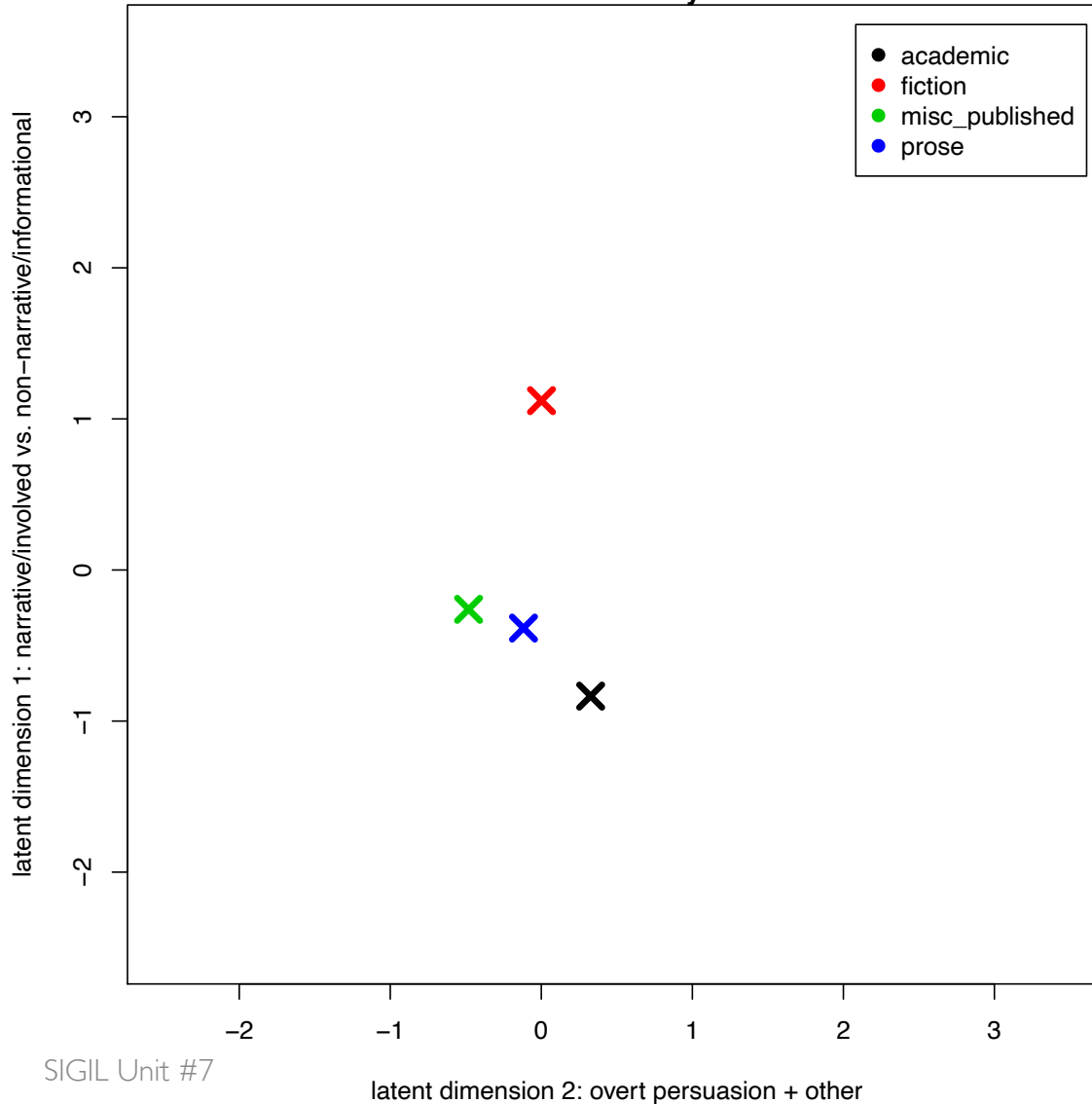


correlated with verb frequency  
(all feat's measured per 1000 words)

- 126 past participle
- 143 type token
- 134 sentence relatives
- 136 though
- 139 prepositions
- 144 mean word length
- 140 adj attr
- 127 past participle whiz
- 118 by passives
- 117 agentless passives
- 164 phrasal coordination
- 114 nominalization
- 145 conjuncts
- 116 other nouns
- 138 other adv sub
- 131 wh subj
- 132 wh obj
- 133 pied piping
- 151 demonstratives
- 157 verb suasive
- 122 that adj comp
- 130 that obj
- 121 that verb comp
- 104 place adverbials
- 105 time adverbials
- 125 present participle
- 147 hedges
- 101 past tense
- 102 perfect aspect
- 108 third person pronouns
- 161 stranded preposition
- 113 wh question
- 107 second person pronouns
- 123 wh clause
- 142 adverbs
- 150 discourse particles
- 159 contractions
- 106 first person pronouns
- 112 proverb do
- 111 indefinite pronoun
- 109 pronoun it
- 167 neg analytic
- 156 verb private
- 149 emphatics
- 155 vto public
- 152 verb seem
- 160 neg synthetic
- 108 present participle whiz
- 114 gerunds
- 144 downtoners
- 148 amplifiers
- 162 split infinitive
- 129 that subj
- 120 existential there
- 135 because
- 103 present tense
- 153 modal necessity
- 152 modal possibility
- 124 infinitives
- 163 split auxiliary
- 154 modal predictive
- 169 if
- 110 demonstrative pronoun
- 119 be main verb
- 141 adj pred

# Design bias: choice of text samples

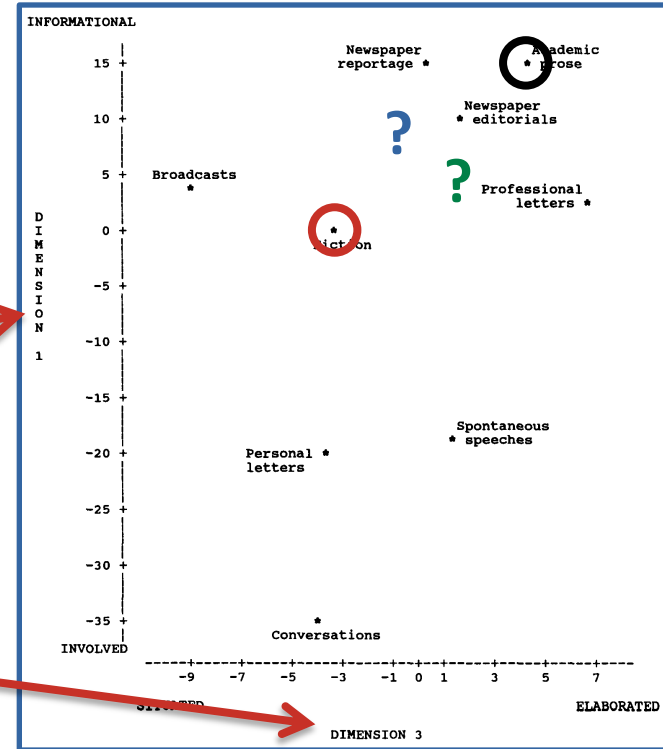
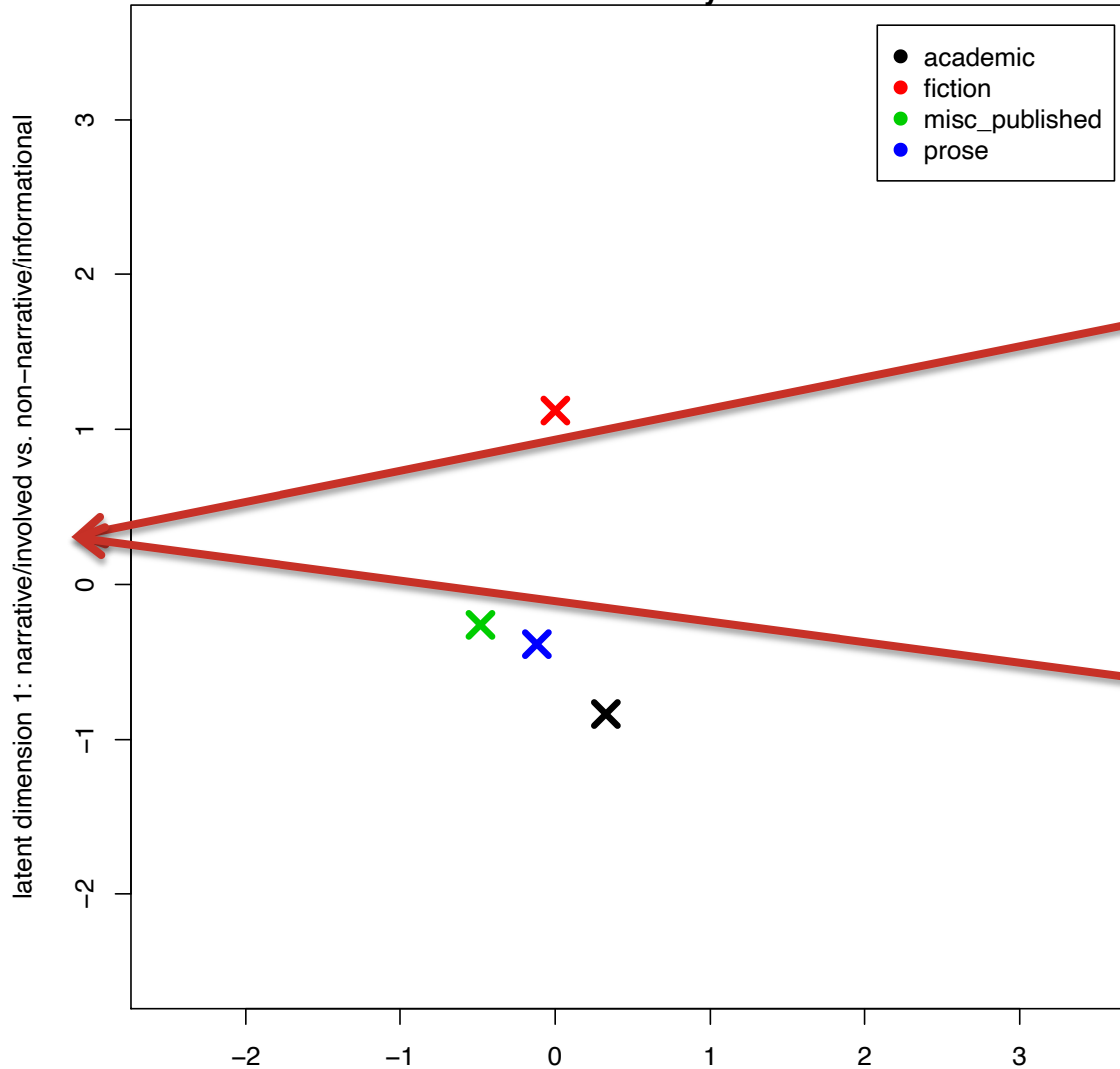
4-Factor Analysis





# Interpretation bias

4-Factor Analysis



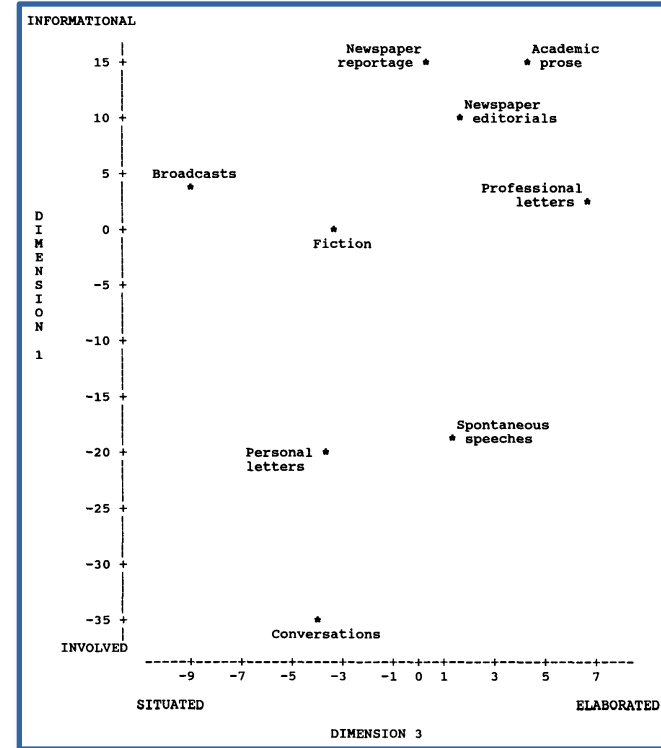
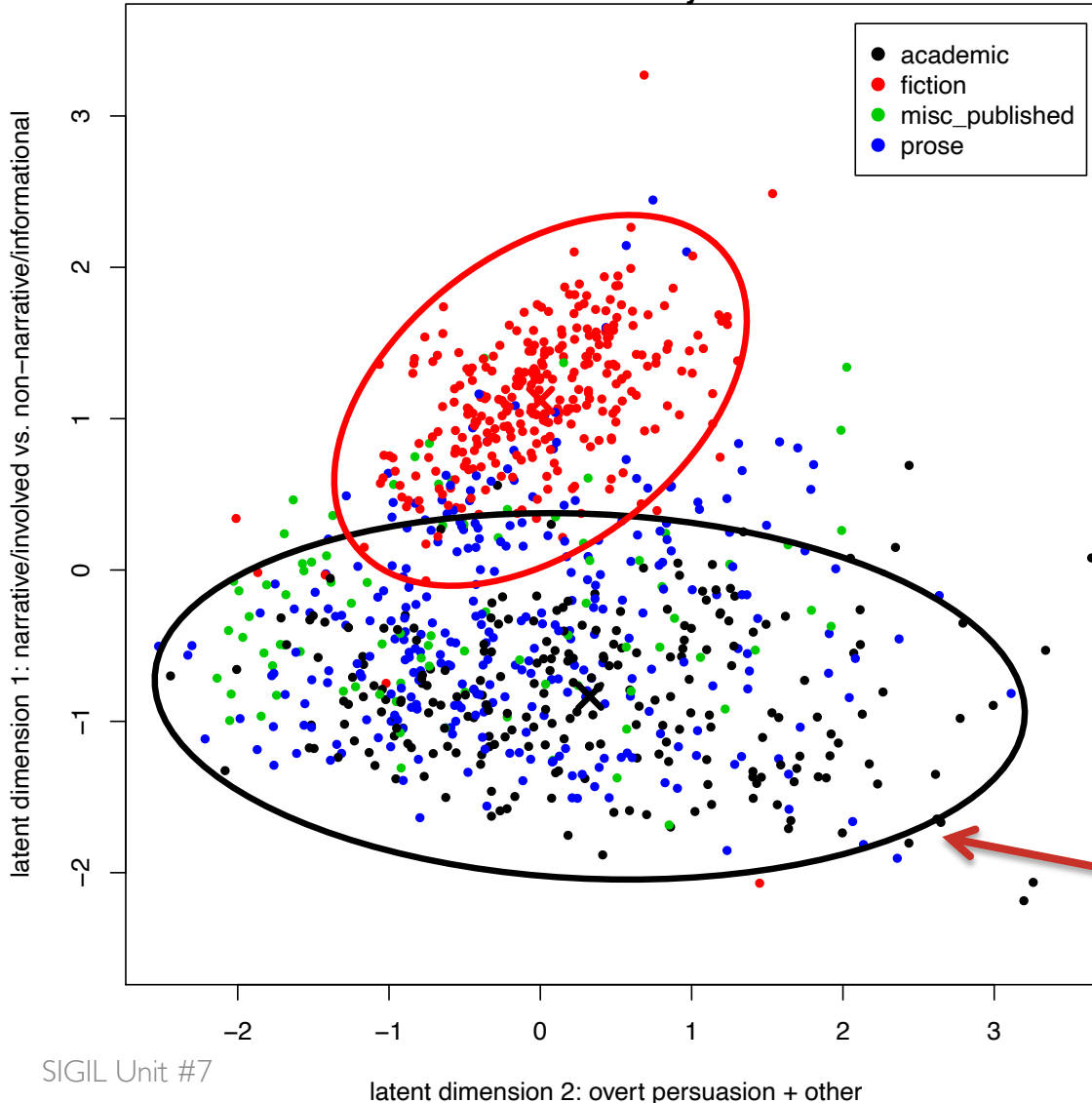
# Interpretation bias

TABLE 2  
 Summary of the co-occurrence patterns underlying five major dimensions of English.

DIMENSION 1 (Informational vs. Involved)		DIMENSION 2 (Narrative versus Non-Narrative)		DIMENSION 3 (Elaborated vs. Situated Reference)		DIMENSION 4 (Overt Expression of Persuasion)		DIMENSION 5 (Absorptive vs. Non-Absorptive)
nouns	0.80	past tense verbs	0.90	WH relative clauses on object positions	0.63	infinitives	0.76	conjunctives
word length	0.58	third person pronouns	0.73	pied piping constructions	0.61	prediction modals	0.54	agentless passives
prepositional phrases	0.54	perfect aspect verbs	0.48	WH relative clauses on subject position	0.45	suasive verbs	0.49	past participles
type / token ratio	0.54	public verbs	0.43	phrasal coordination	0.36	conditional subordination	0.47	clausal adjuncts
attributive adjs.	0.47	synthetic negation	0.40	nominalizations	0.36	necessity modals	0.46	BY-passives
		present participial clauses	0.39	time adverbials	-0.60	split auxiliaries	0.44	past participles with WH
private verbs	-0.96	present tense verbs	-0.47	place adverbials	-0.49	possibility modals	0.37	other auxiliaries
<i>that</i> deletions	-0.91	attributive adjs.	-0.41	other adverbs	-0.46	[No complementary features]		subordinate clauses
contractions	-0.90							
present tense verbs	-0.86							
2nd person pronouns	-0.86							
<i>do</i> as pro-verb	-0.82							
analytic negation	-0.78							
demonstrative pronouns	-0.76							
general emphatics	-0.74							
first person pronouns	-0.74							
pronoun <i>it</i>	-0.71							
<i>be</i> as main verb	-0.71							
causative subordination	-0.66							
discourse particles	-0.66							
indefinite pronouns	-0.62							
general hedges	0.58							

# Variation between texts is ignored

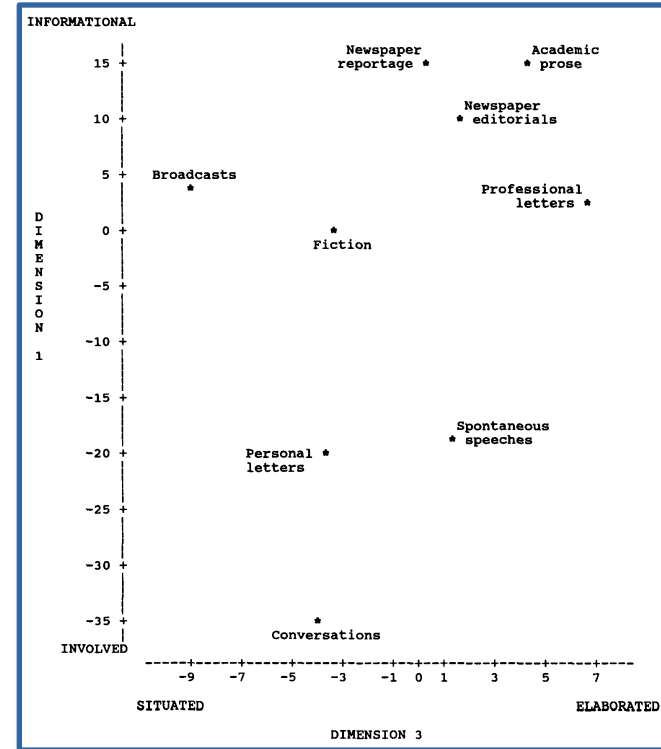
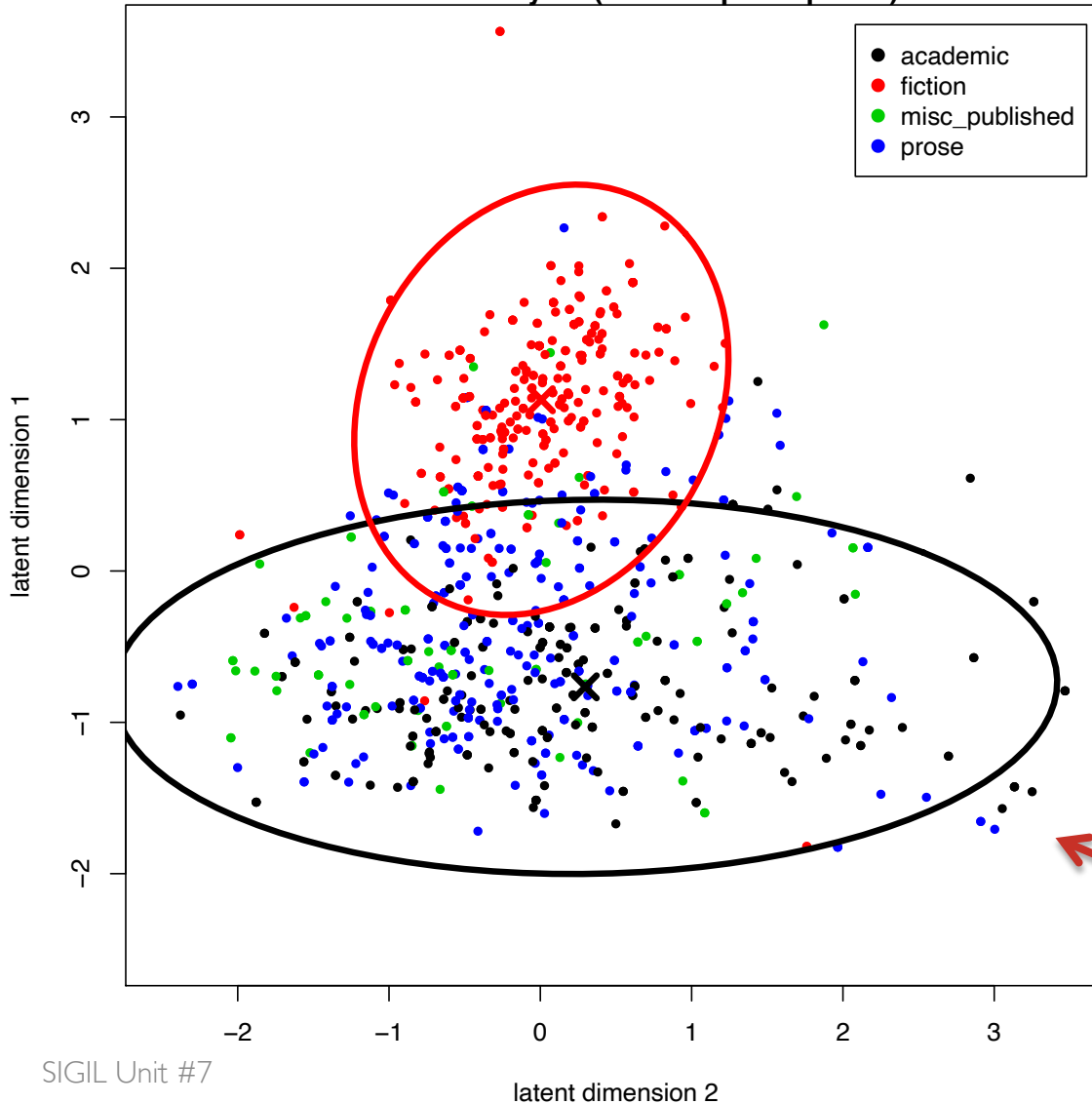
4-Factor Analysis



“confidence” ellipse  
(→ significance)

# Design bias: choice of texts (redux)

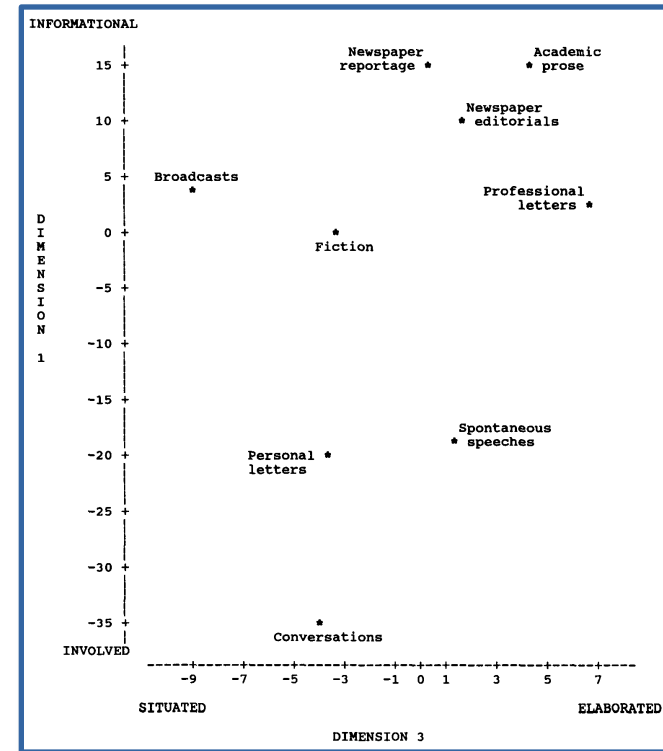
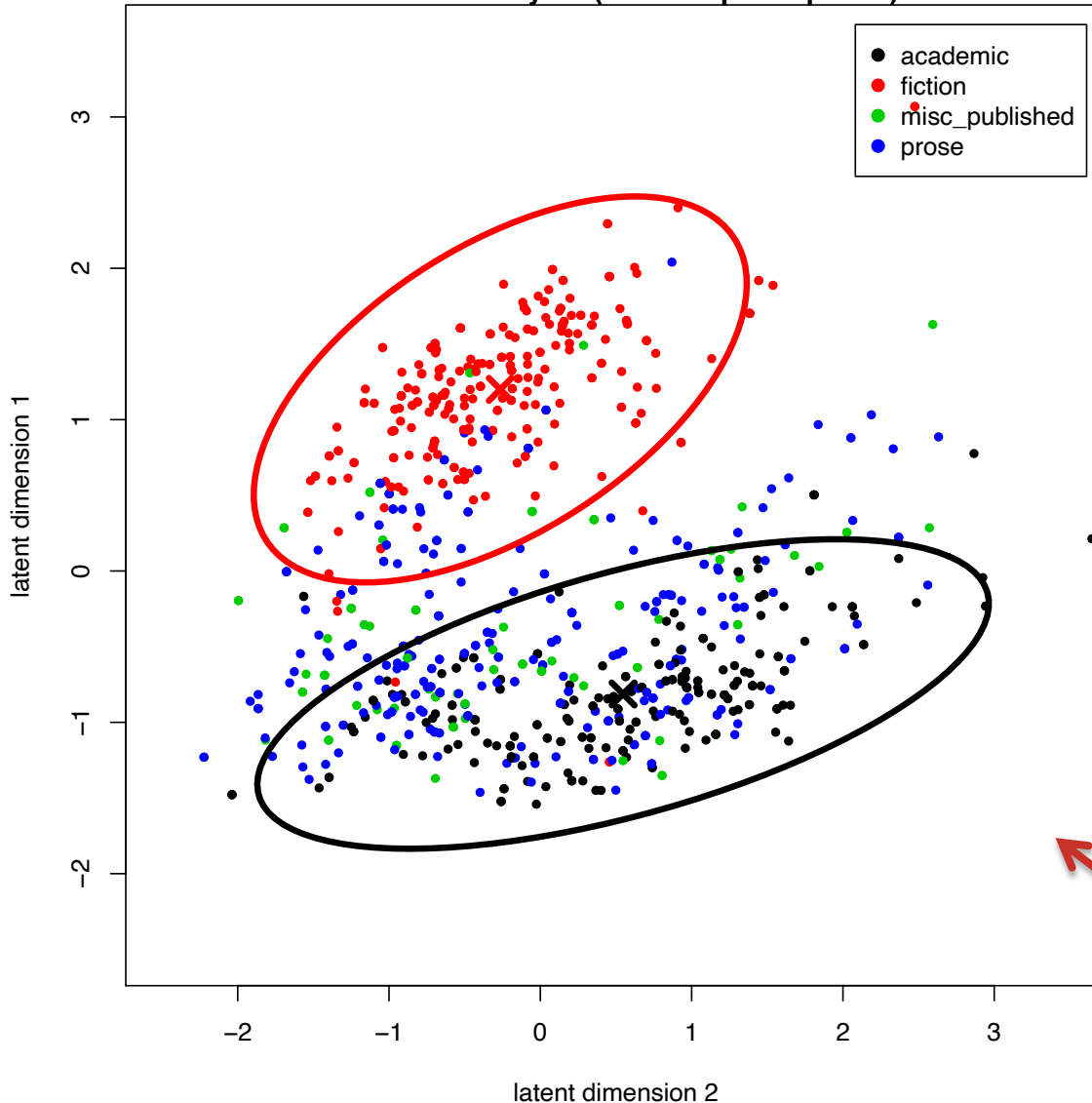
4-Factor Analysis (bootstrap sample #3)



Bootstrapping

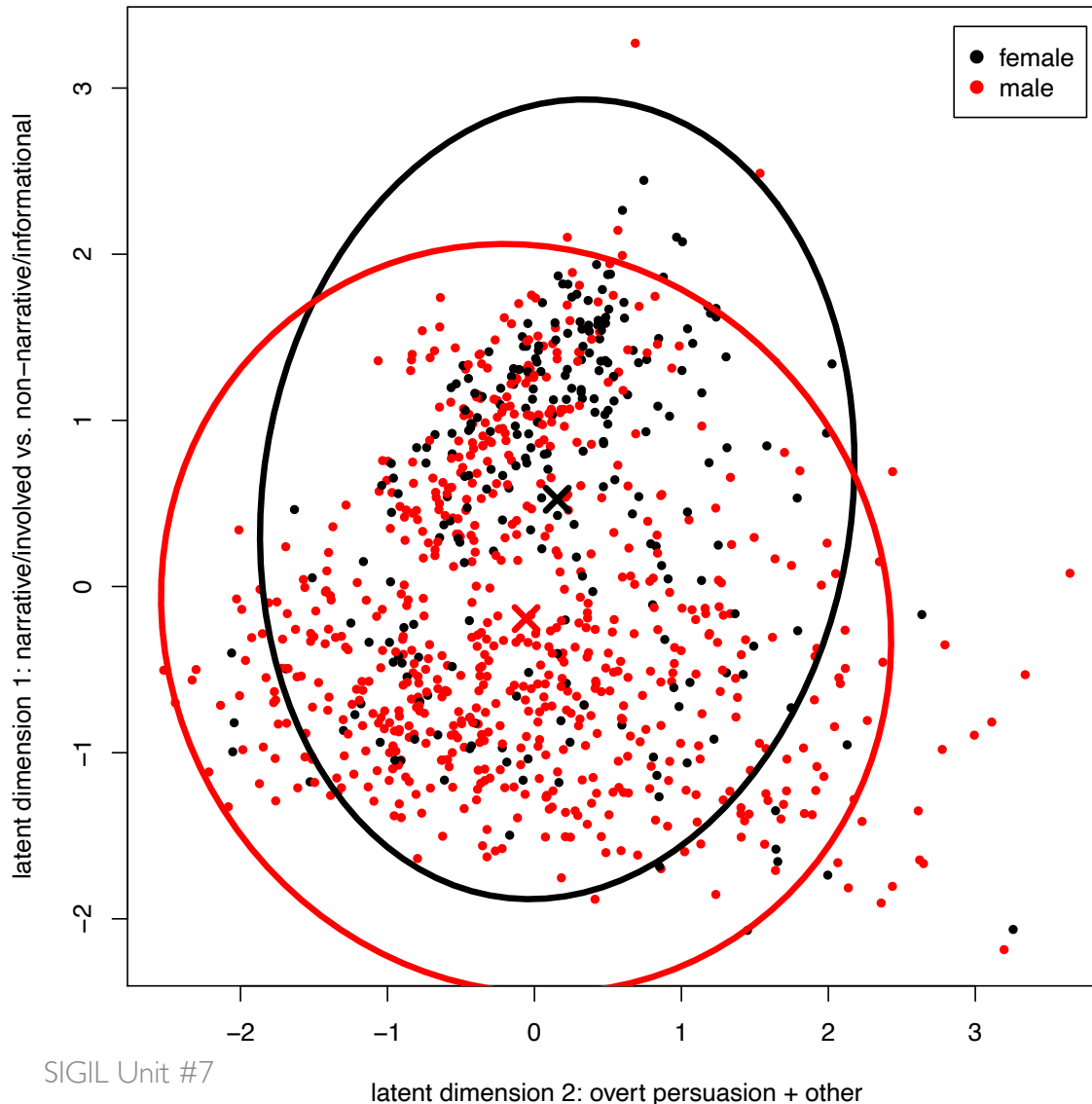
# And there's the magic number ...

3-Factor Analysis (bootstrap sample #3)



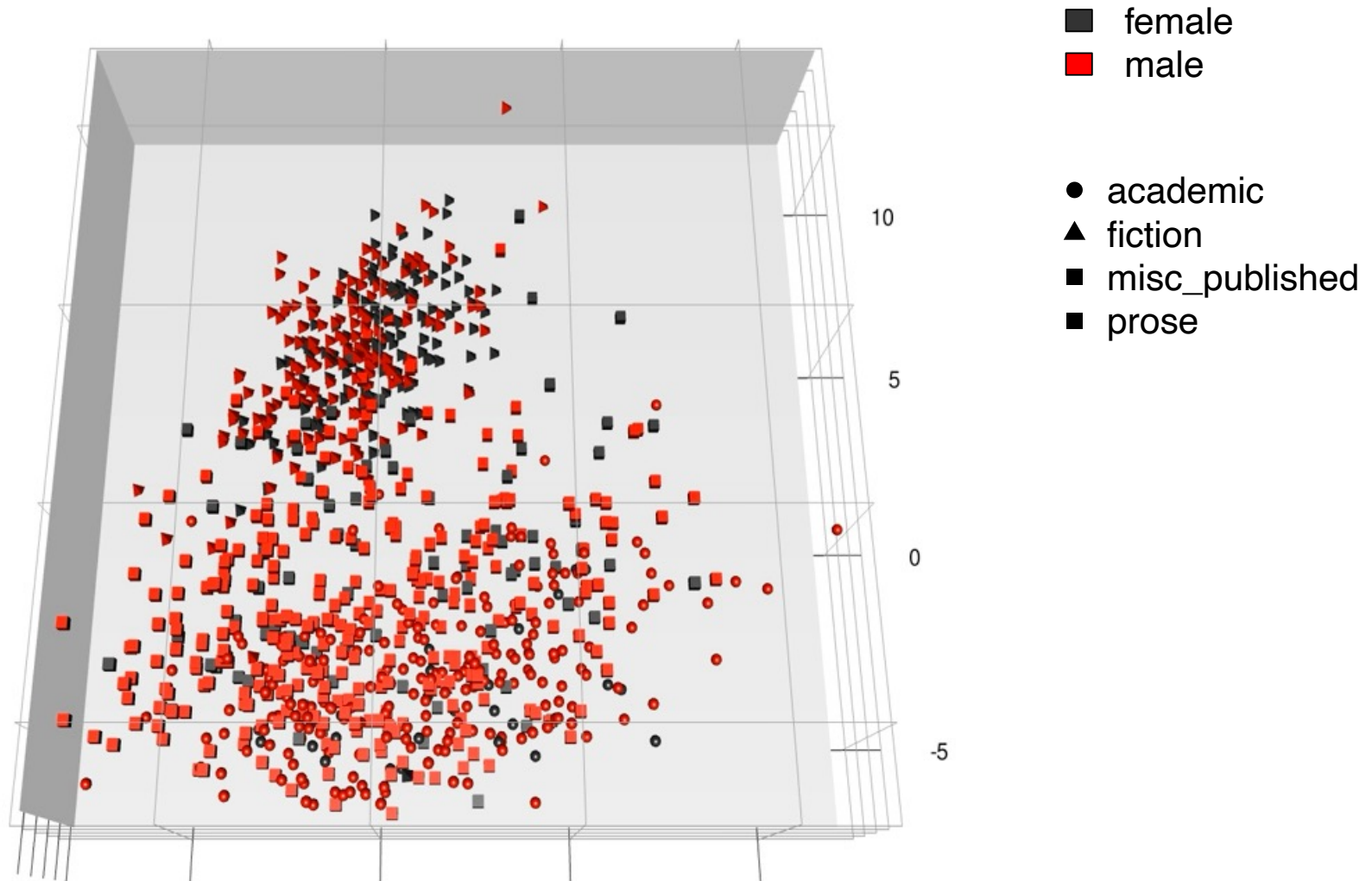
3-factor analysis  
(instead of 4 factors)

# Blindness to subtle patterns



- But research shows that author gender can be identified with high accuracy
  - Koppel et al. (2003): 77.3% with function words + POS n-grams
  - Gasthaus (2007): 82.9% with SVM on Biber features
- This dataset:
  - 82.3% accuracy
  - baseline: 73.1%

# Blindness to subtle patterns



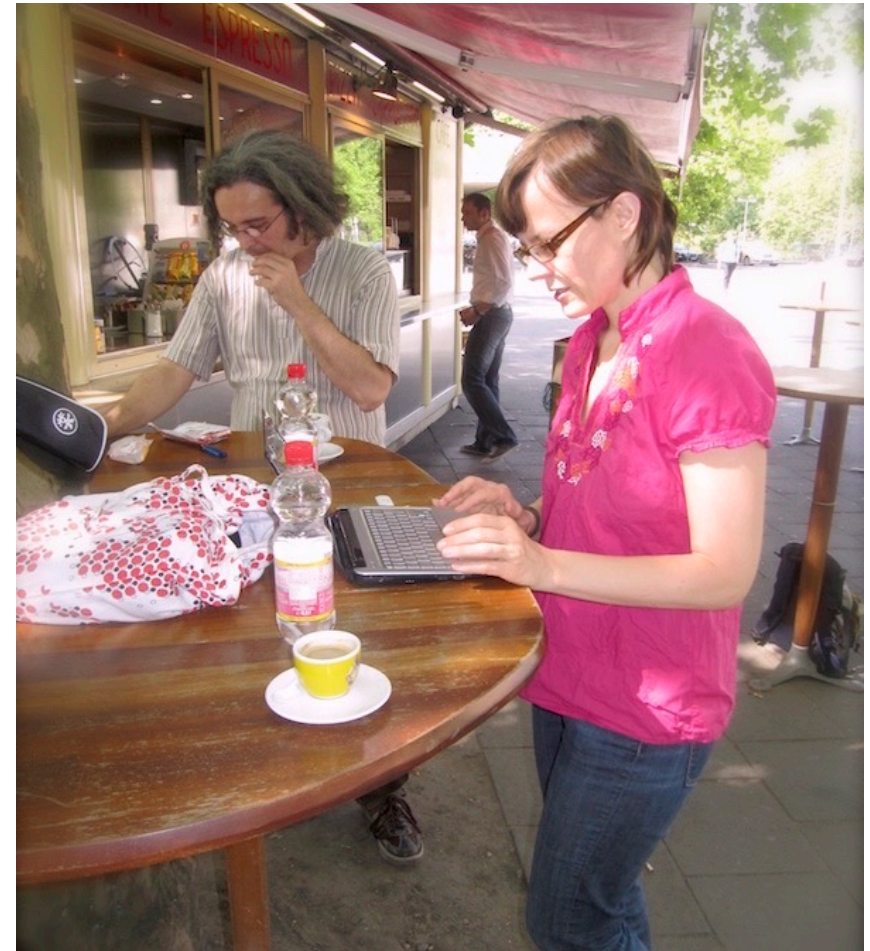
# Geometric Multivariate Analysis

(Diwersy, Evert & Neumann 2014; Evert & Neumann 2017; Neumann & Evert 2021)

Online supplements:

[https://www.stephanie-evert.de/  
PUB/EvertNeumann2017/](https://www.stephanie-evert.de/PUB/EvertNeumann2017/)

[https://www.stephanie-evert.de/  
PUB/NeumannEvert2021/](https://www.stephanie-evert.de/PUB/NeumannEvert2021/)





# Geometric Multivariate Analysis

(Diwersy, Evert & Neumann 2014; Evert & Neumann 2017; Neumann & Evert 2021)

- Axiom: (Euclidean) distance = similarity of texts
  - depends crucially on theoretically motivated features
- Visualization → interpret geometric configuration
  - latent dimensions as geometric projections
  - orthogonal projection = perspective on data
  - method: principal component analysis (PCA)
- Minimally supervised intervention
  - based on externally observable, theory-neutral information
  - method: linear discriminant analysis (LDA)
- Bootstrapping / cross-validation to assess significance
- Cautious interpretation of feature weights

# Case study: CroCo

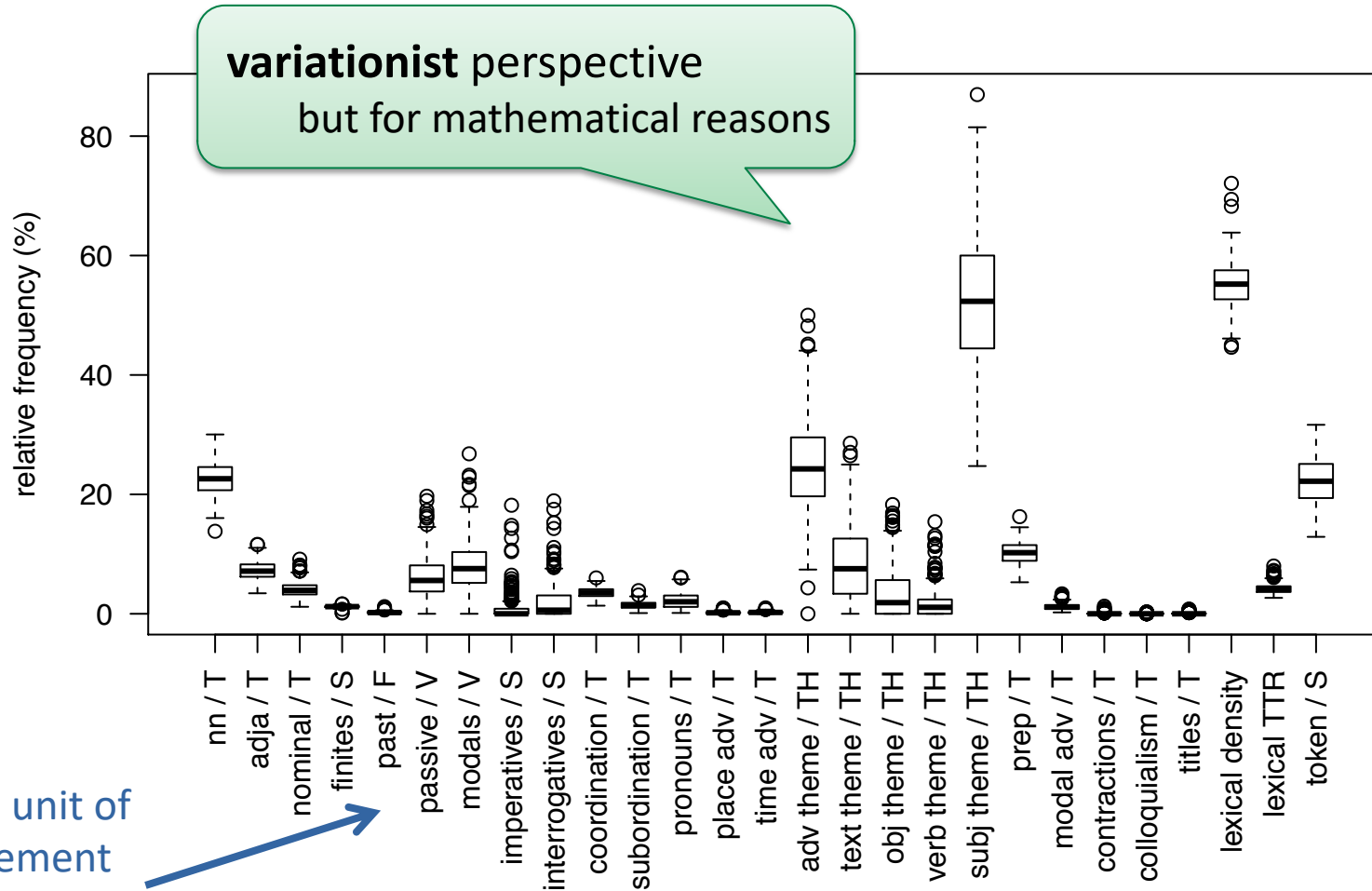
(Neumann 2013; Evert & Neumann 2017)

- CroCo: parallel corpus English/German
  - English-German and German-English translation pairs
  - we use 298 texts from 5 different genres (excluded: instruction manuals, tourism, fiction)
- 28 lexico-grammatical features (Neumann 2013)
  - comparable between languages
  - inspired by SFL and translation studies
- Text = point in 28-dimensional feature space
- Research hypotheses: **shining through** (Teich 2003), **prestige effect** (Toury 2012)

**genre:** language-external  
situation + purpose

**register:** language-internal  
co-occurrence patterns of  
linguistic features

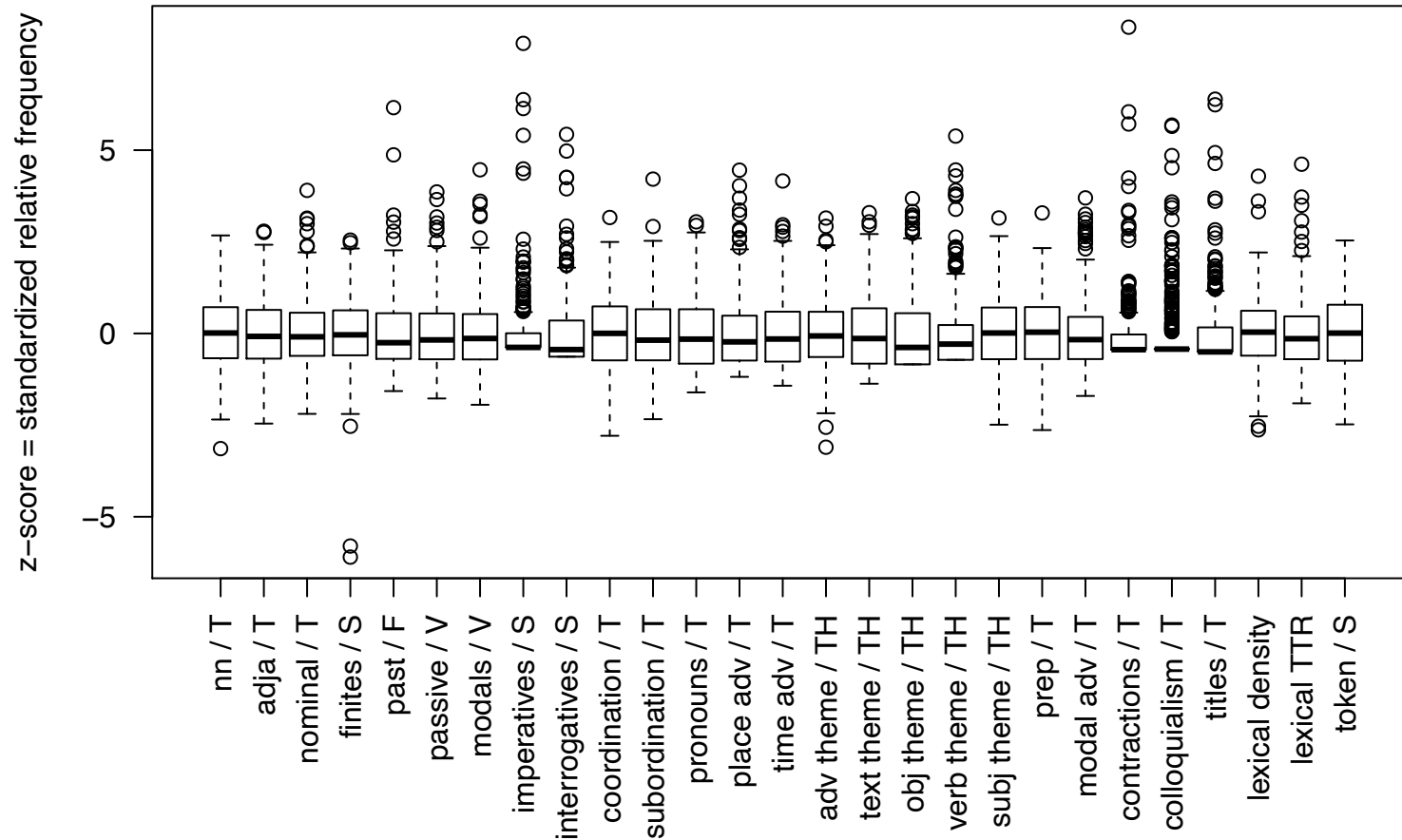
# Feature design: avoid “obvious” correlations



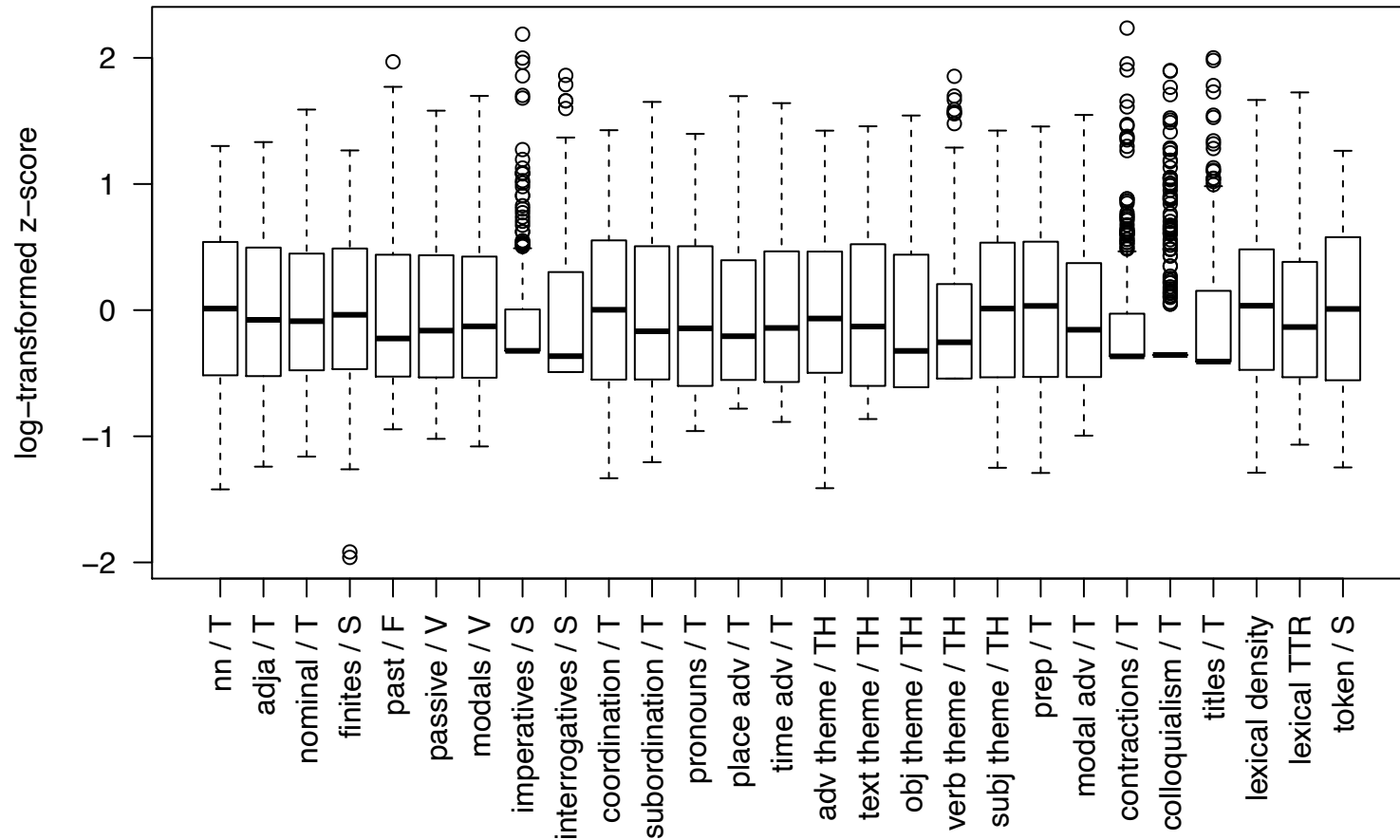
suitable unit of  
measurement  
(not always per  
1000 words!)

# Feature scaling:

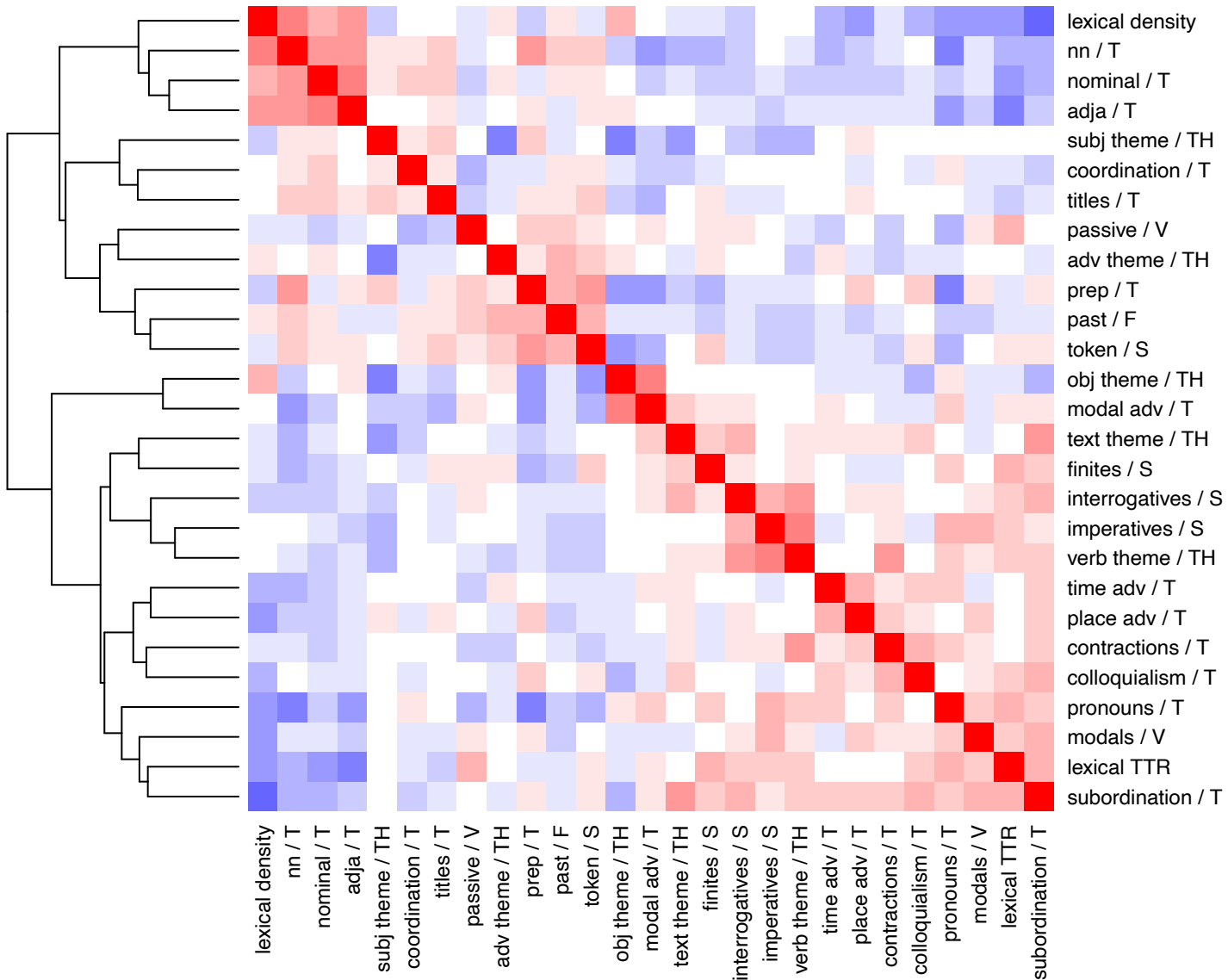
## same contribution to Euclidean distances



# Feature scaling: optional signed log transformation



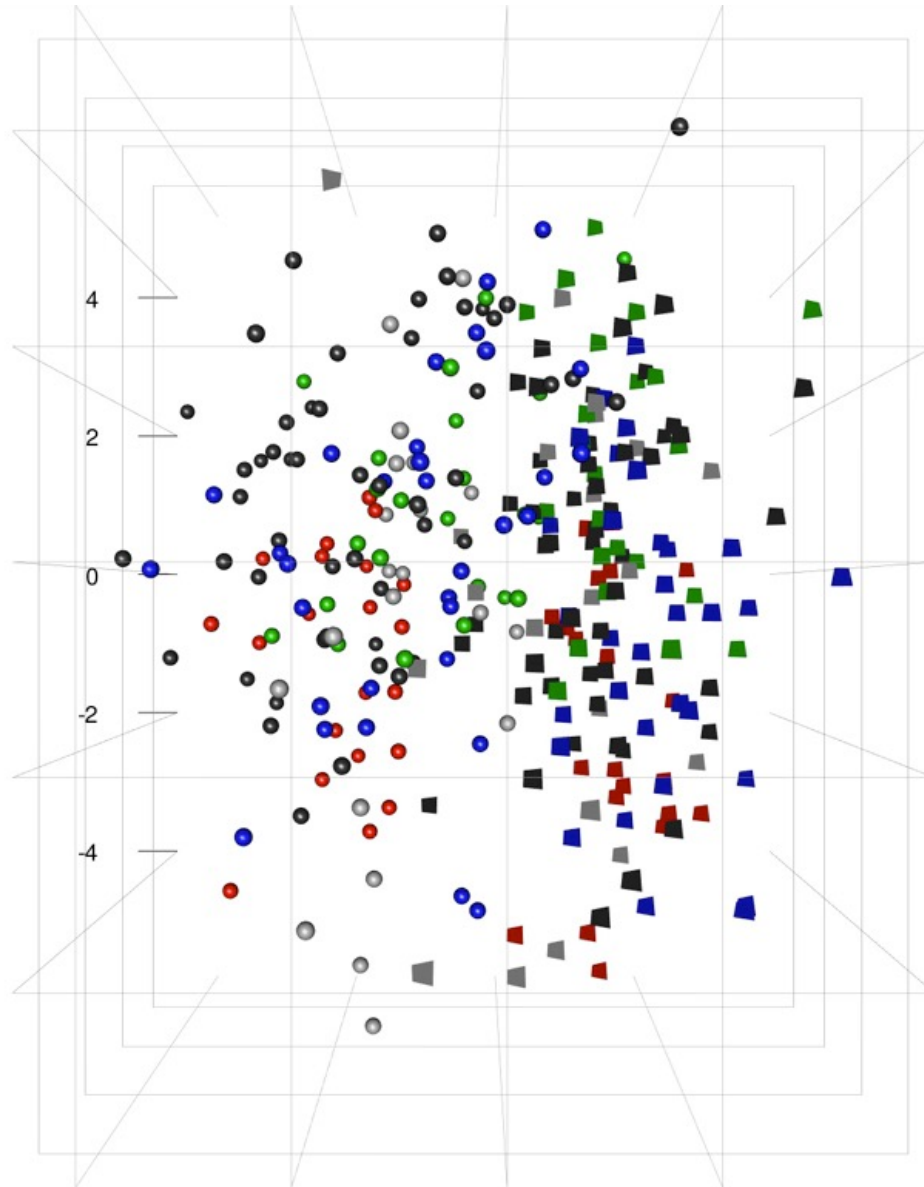
# CroCo: correlation matrix



# Latent dimensions as perspective on data configuration

- Instead of “magical” latent dimensions we focus on **orthogonal projections** as perspectives on the data
  - cf. photograph as 2D perspective on 3D object
- Different perspectives highlight different aspects
- Multivariate analysis → choice of perspective
  - **principal component analysis** (PCA) = perspective that reflects distances between texts as accurately as possible
  - should reveal major dimensions of variation
  - advantage over factor analysis (FA):  
dimensionality does not have to be fixed *a priori*

# CroCo: 3-dimensional projection

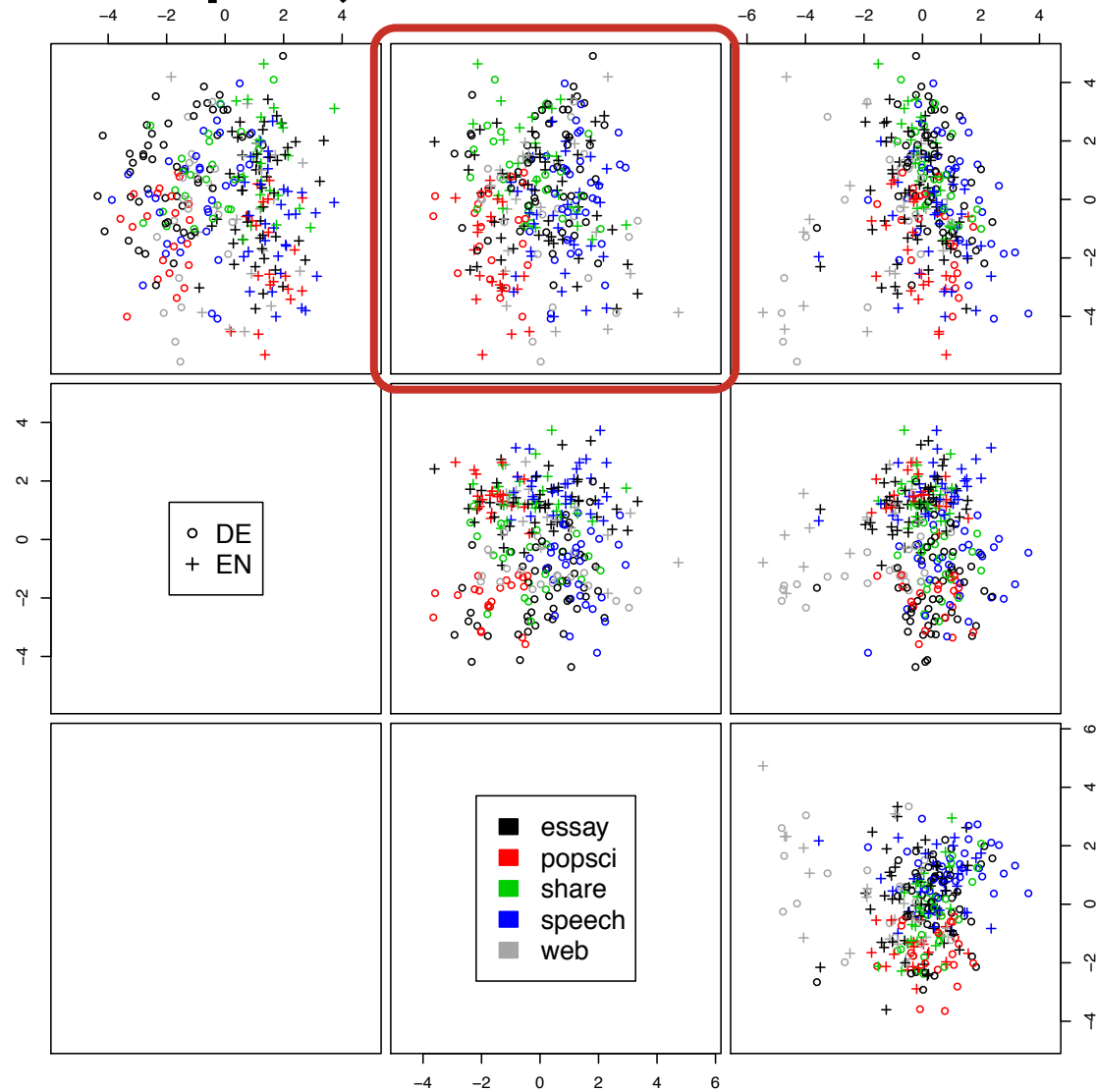
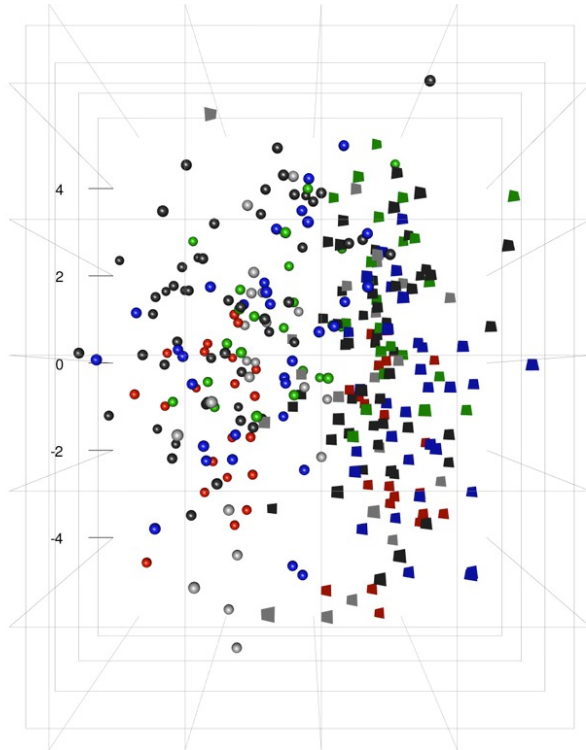


● DE  
▲ EN

■ essay  
■ popsci  
■ share  
■ speech  
■ web

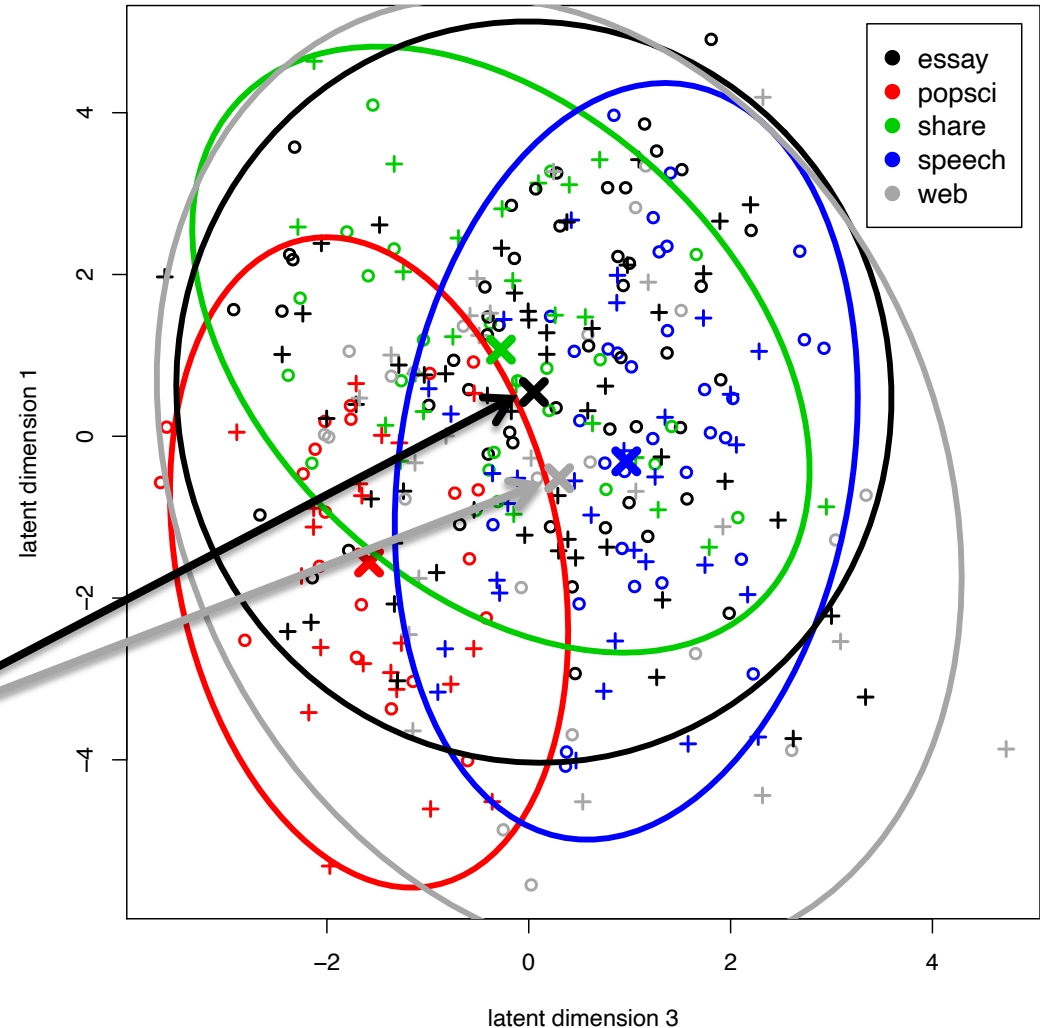


# CroCo: 4-dimensional projection



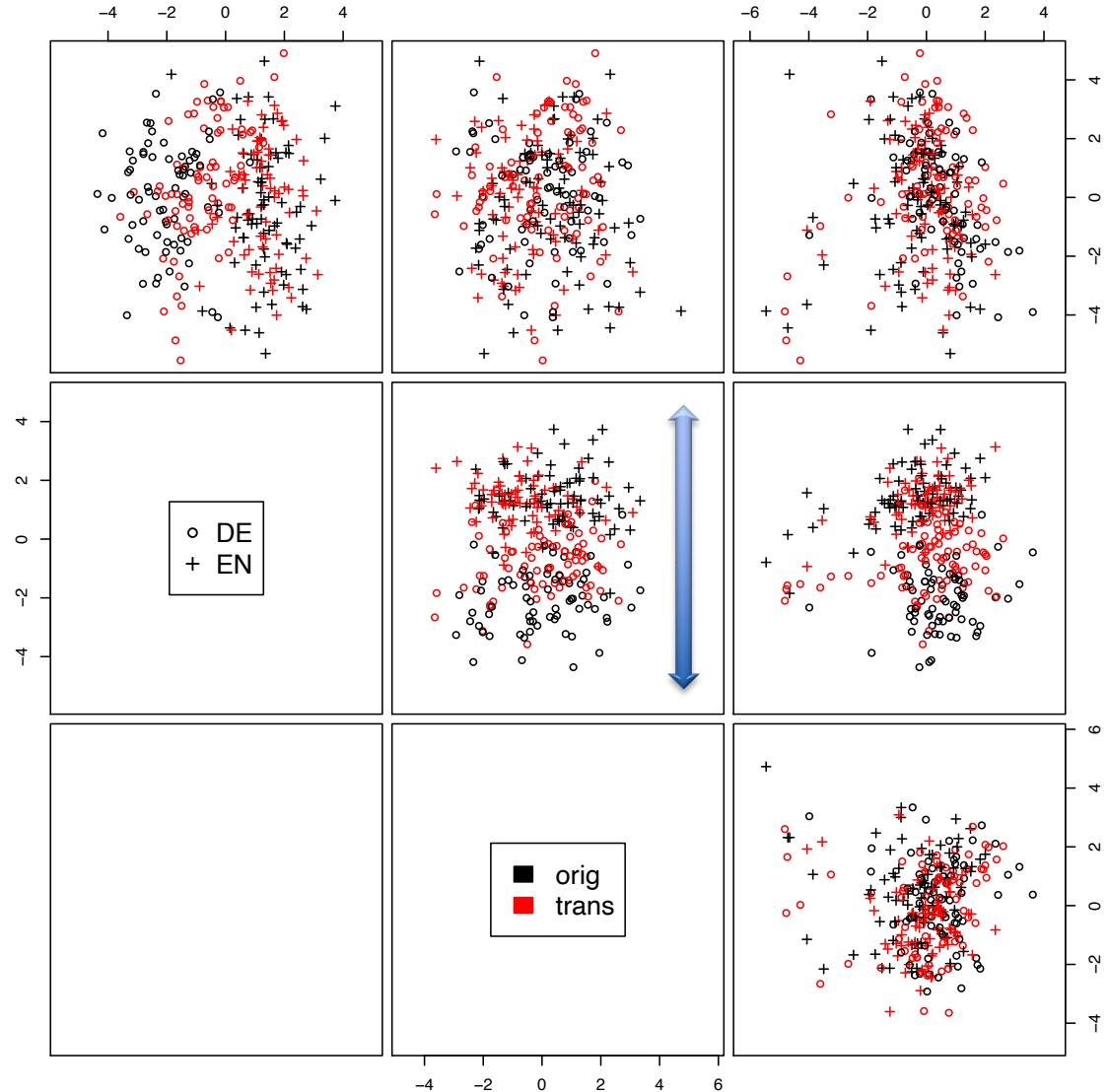
# CroCo: genre distribution

- Focus on latent dim's 1 and 3 (register variation)
- Describe genre by centroid + ellipse
- Comparison with Hotelling's  $t^2$  test
  - essays vs. Web
  - $t^2=4.21$ ,  $df=2/141$ ,  $p=.0167$  \*



# How about translationese?

- PCA dim's can't separate translations from original texts
  - 62.1% accuracy on first 3 PCA dim's
- But SVM machine learner can do this with >80% accuracy
  - RBF kernel
  - 10-fold c.v.
- Hints at **shining through**, but no clear-cut evidence



# Minimally supervised LDA

- Add minimal amount of supervised knowledge to find a more informative perspective
  - evidence for shining through hypothesis from dimension that corresponds to contrast German vs. English
  - supervised knowledge: language of **original texts** only
- Linear **discriminant** analysis (LDA)
  - maximize separation between German / English originals
  - minimize variability within each group
  - classical technique related to PCA and ANOVA
- Project *all* texts onto LDA discriminant
  - complemented by additional PCA dim's for visualization

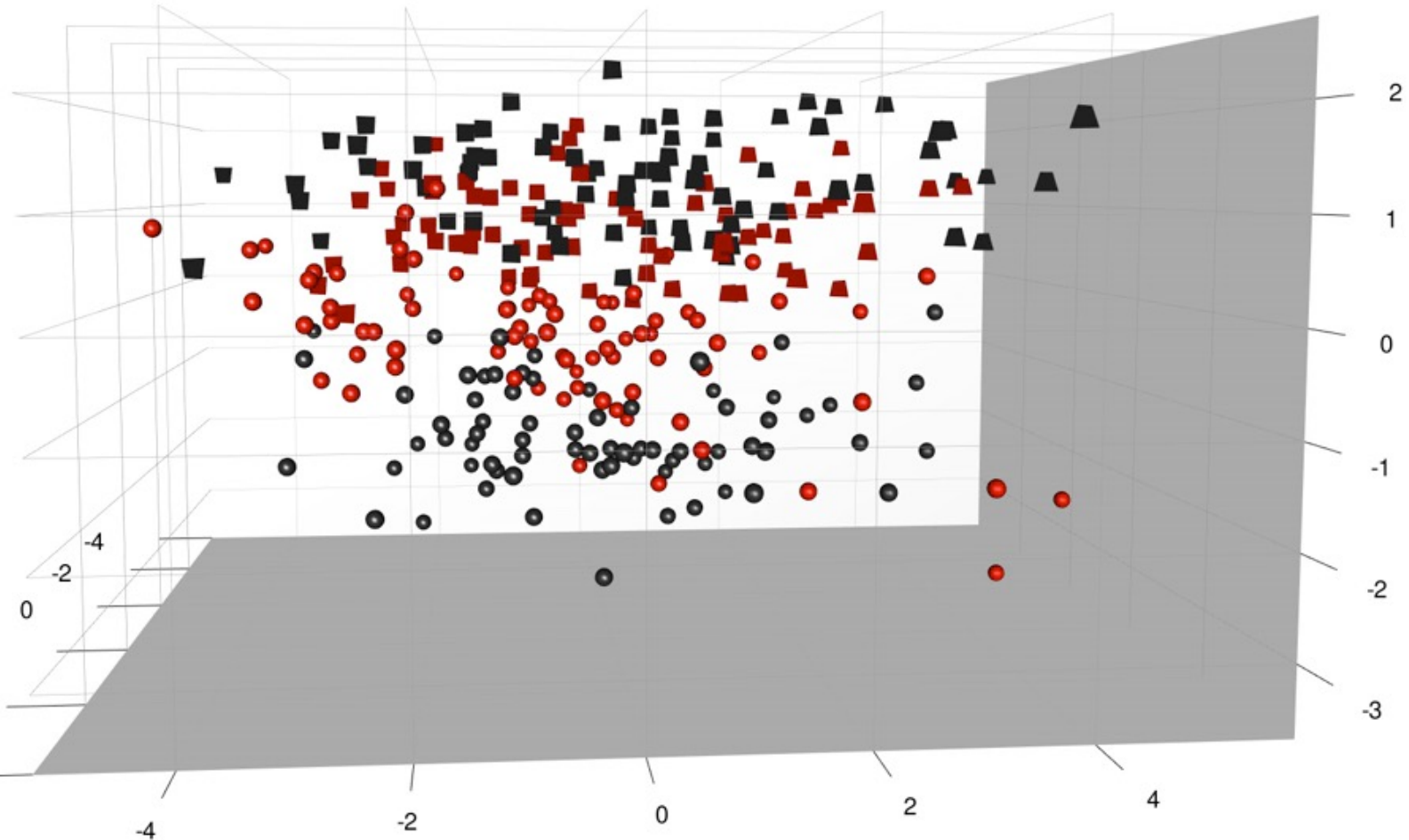
# CroCo: LDA perspective

- DE
- ▲ EN
- orig
- trans

English



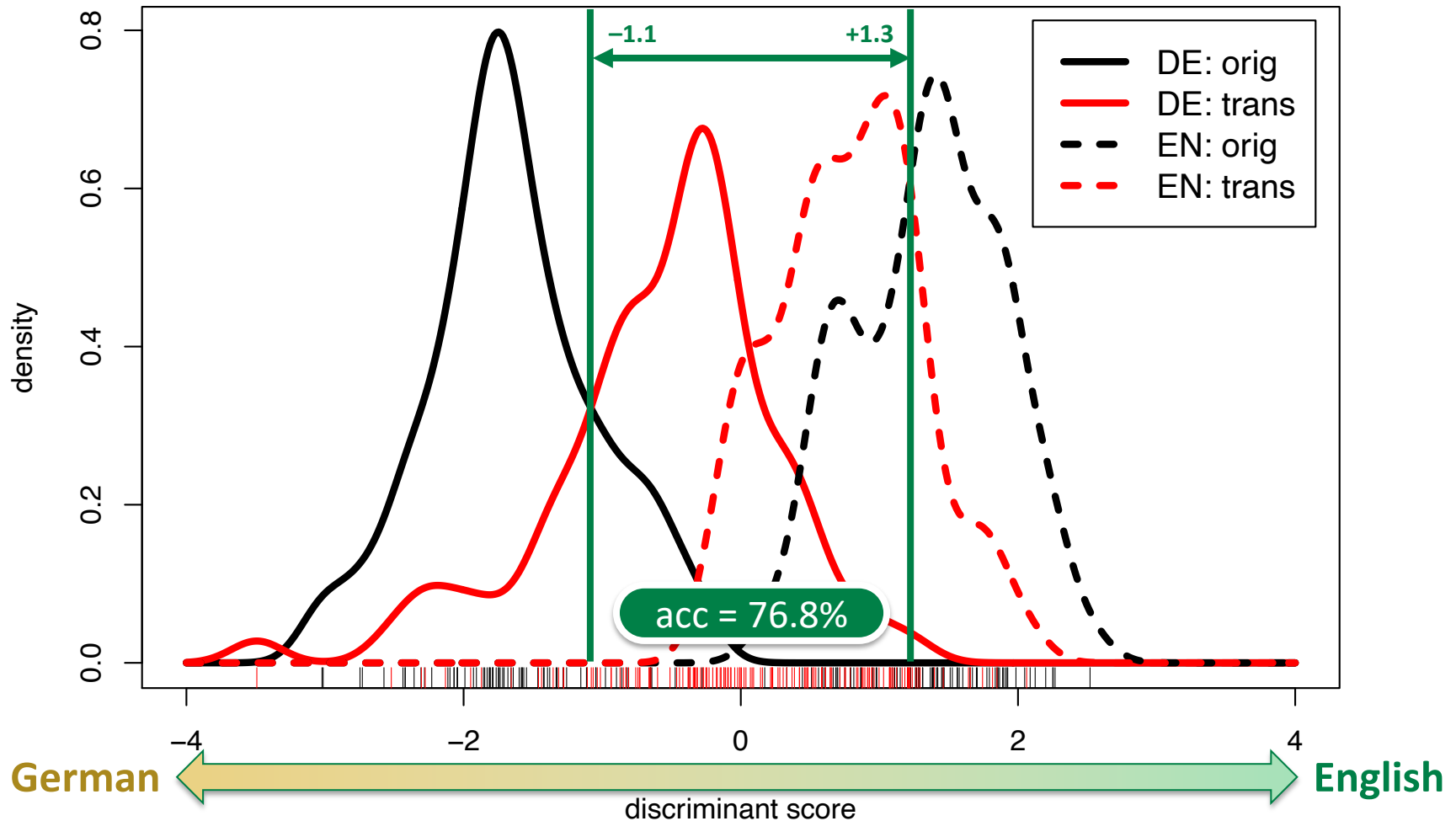
German



# Discriminant for DE vs. EN

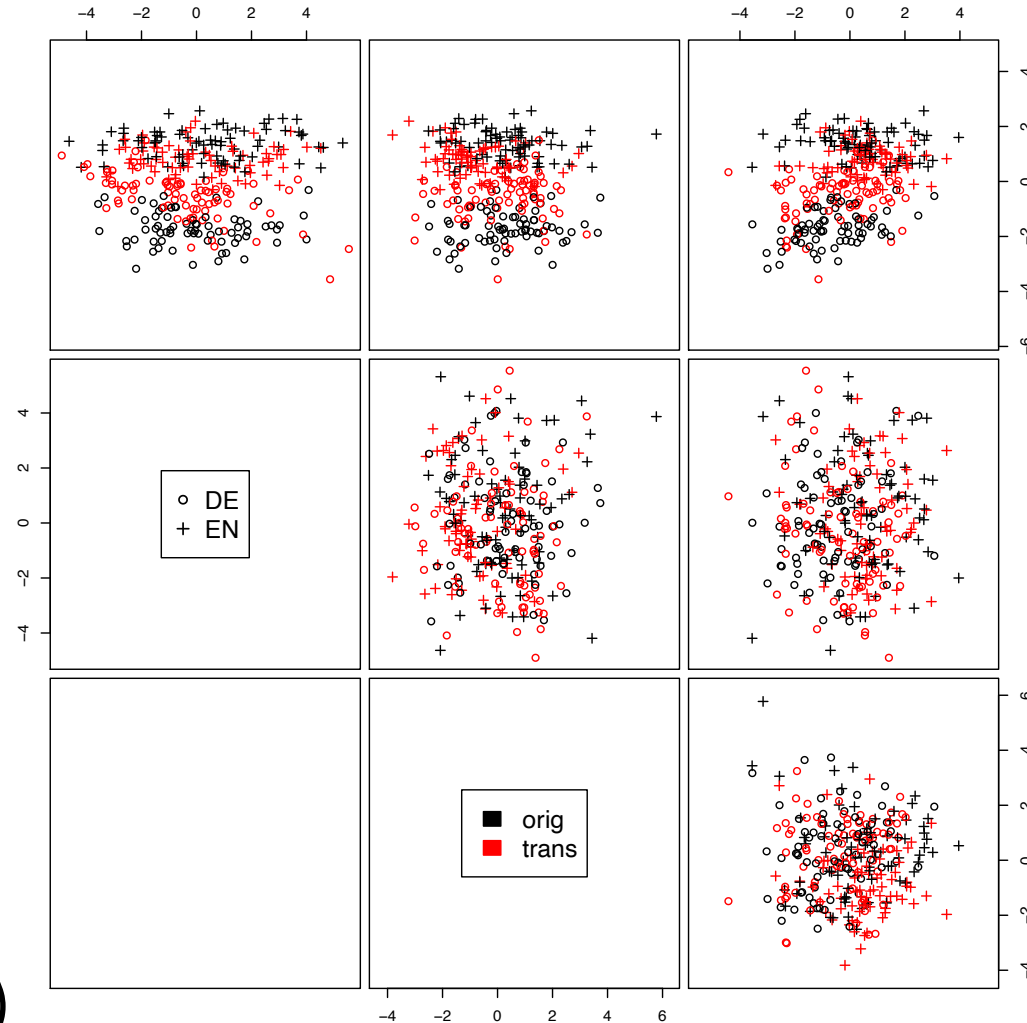


## confirms shining through & prestige effect



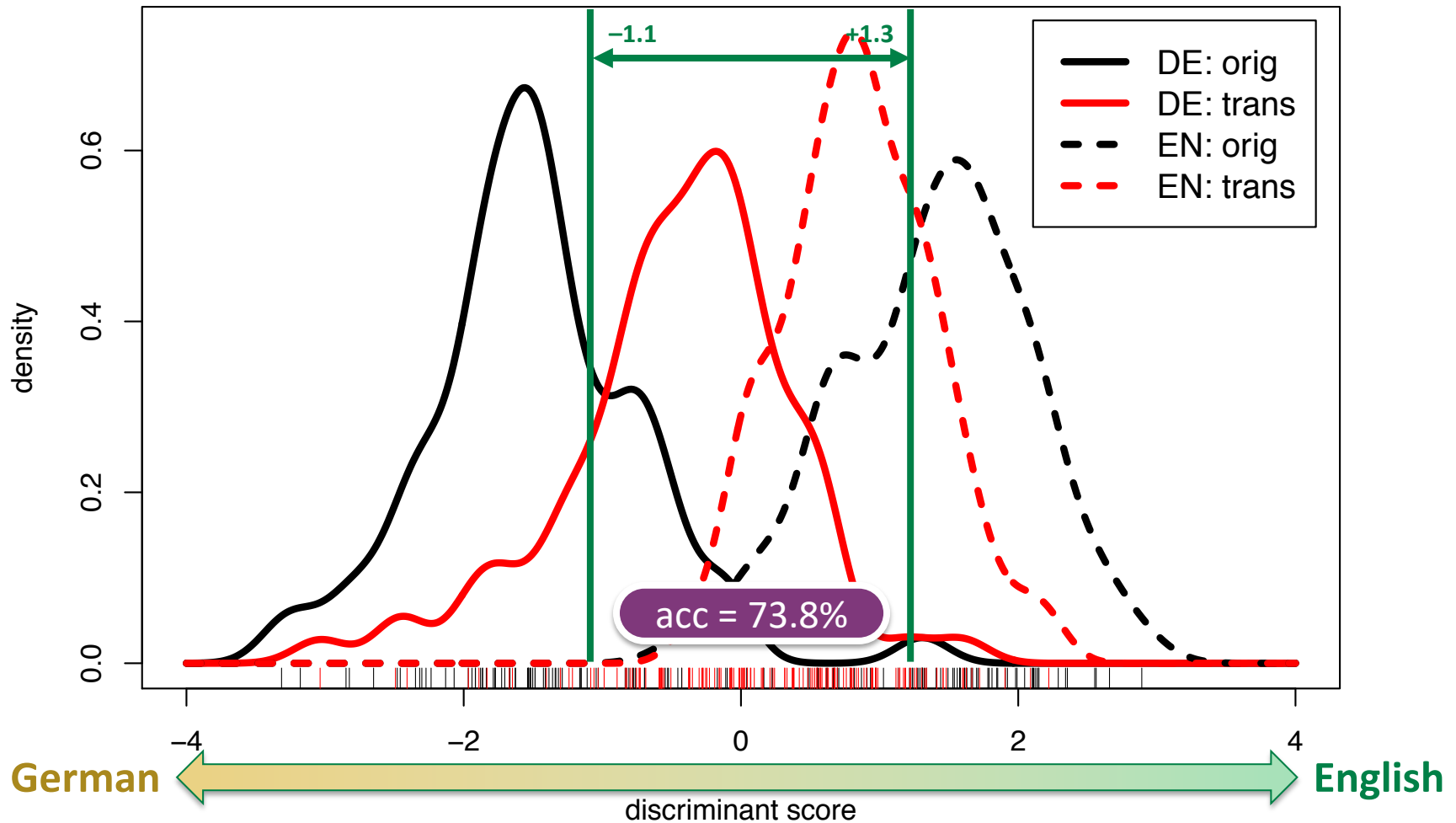
# LDA significance: bootstrapping / cross-validation

- LDA is a supervised ML technique → overtrained?
  - Would we find the same discriminant if we trained on a different set of texts?
- Verification with **bootstrap resampling** or **10-fold cross-validation**
  - LDA trained on 90% of data
  - discriminant axis shows “wobble” of approx. 10°
- Discriminant scores from c.v. (10% test data per fold)



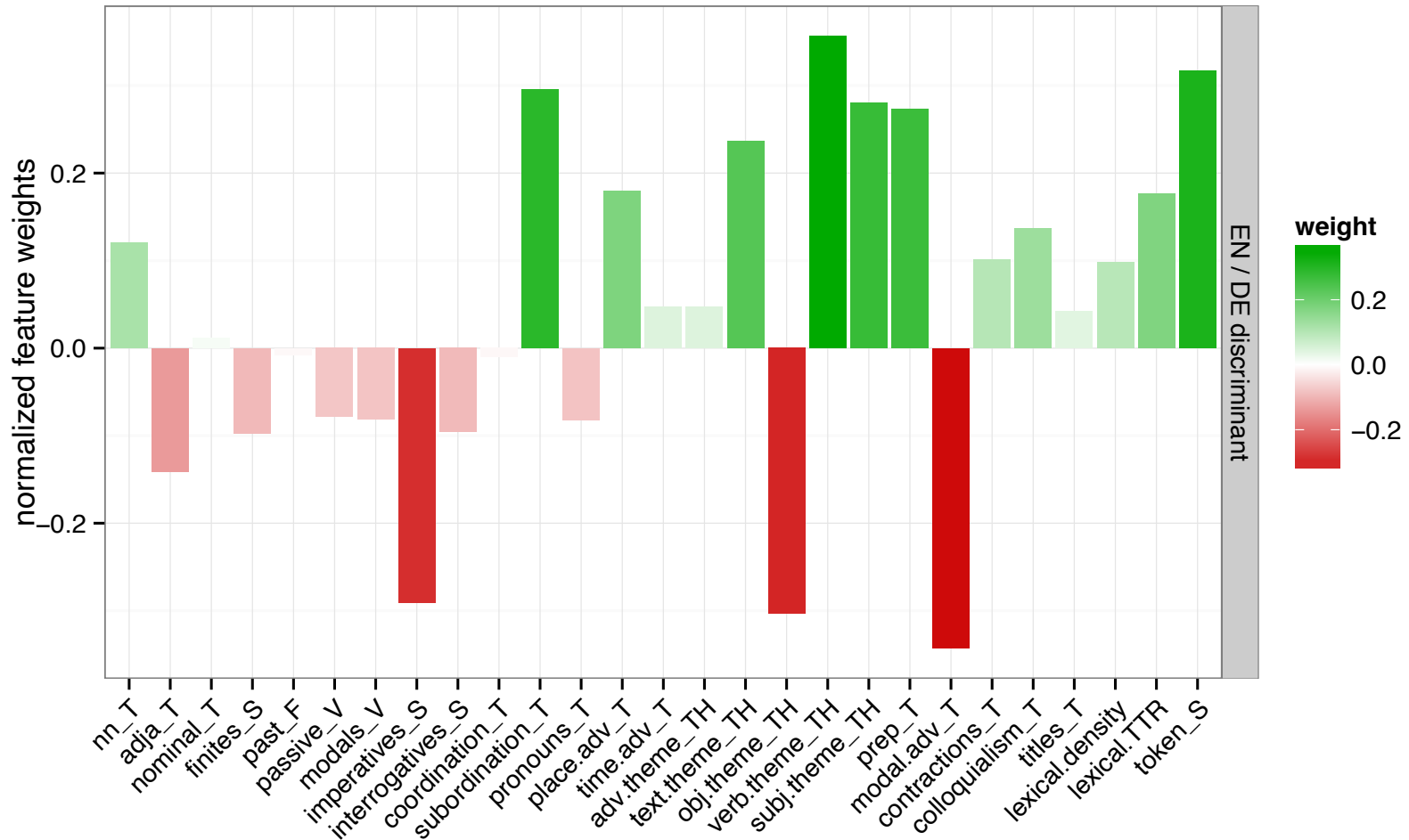
# Cross-validated discriminant

10-fold cross-validation

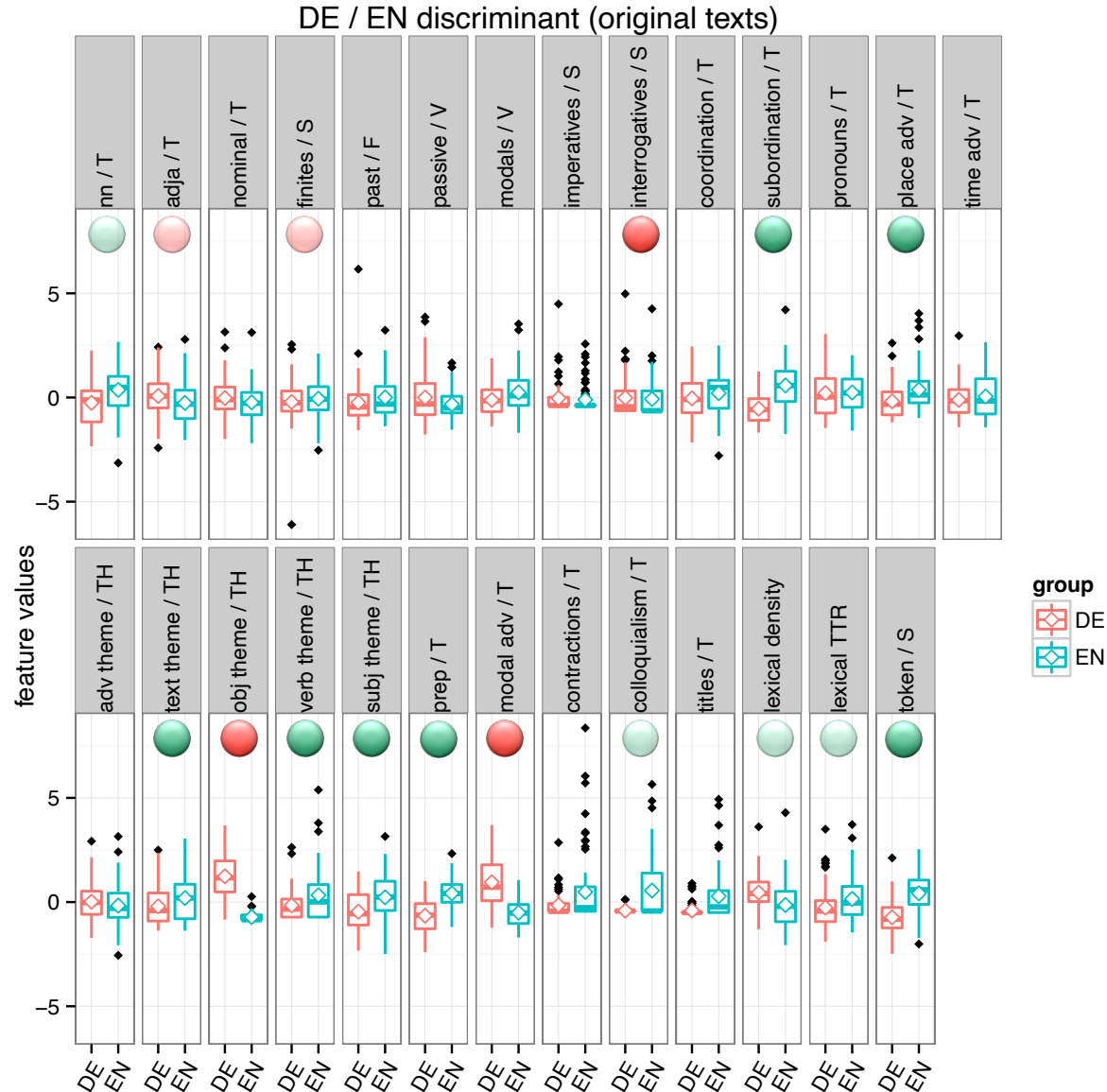




# Interpreting discriminant features

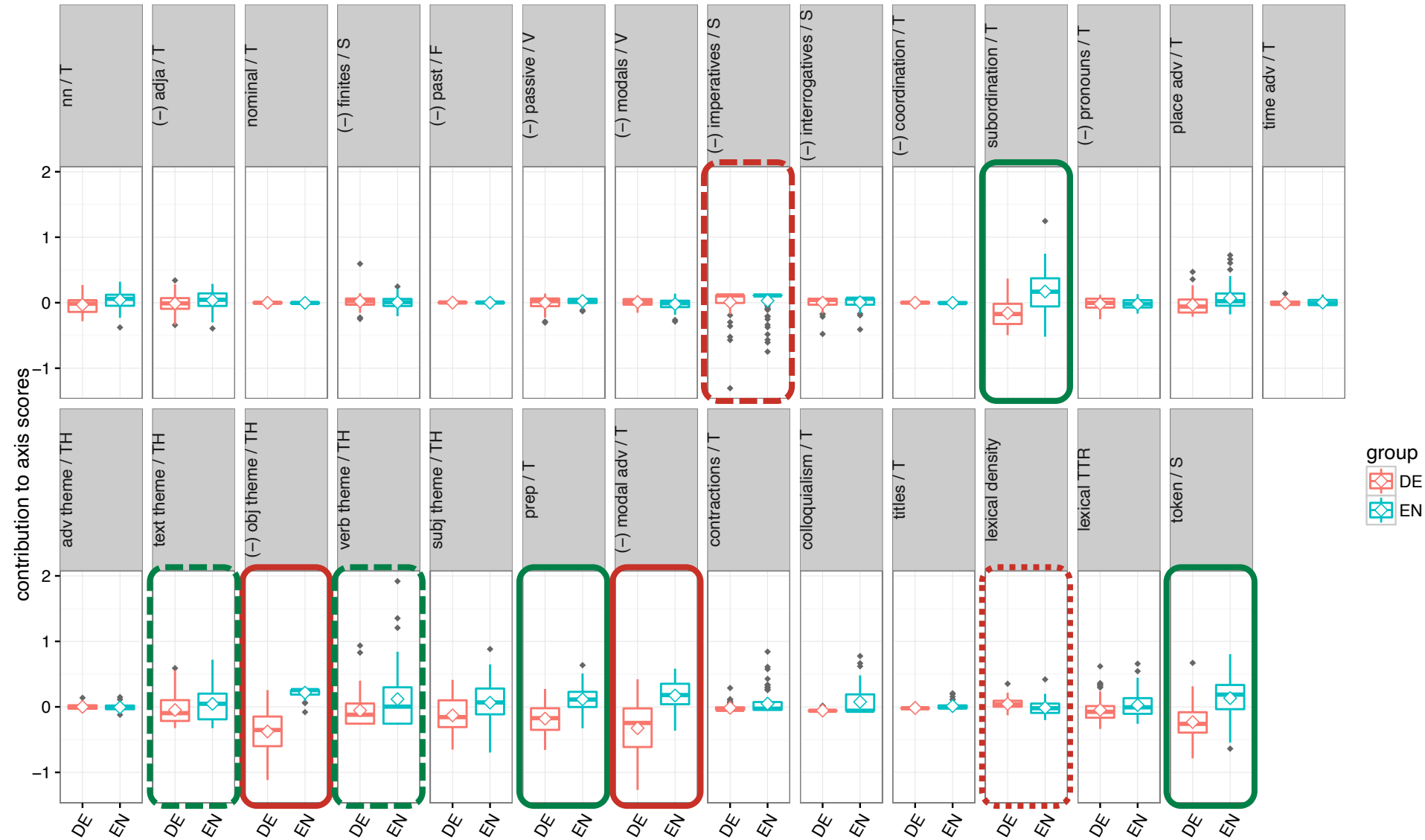


# Interpreting discriminant features



# Interpreting discriminant features

DE / EN discriminant (original texts)



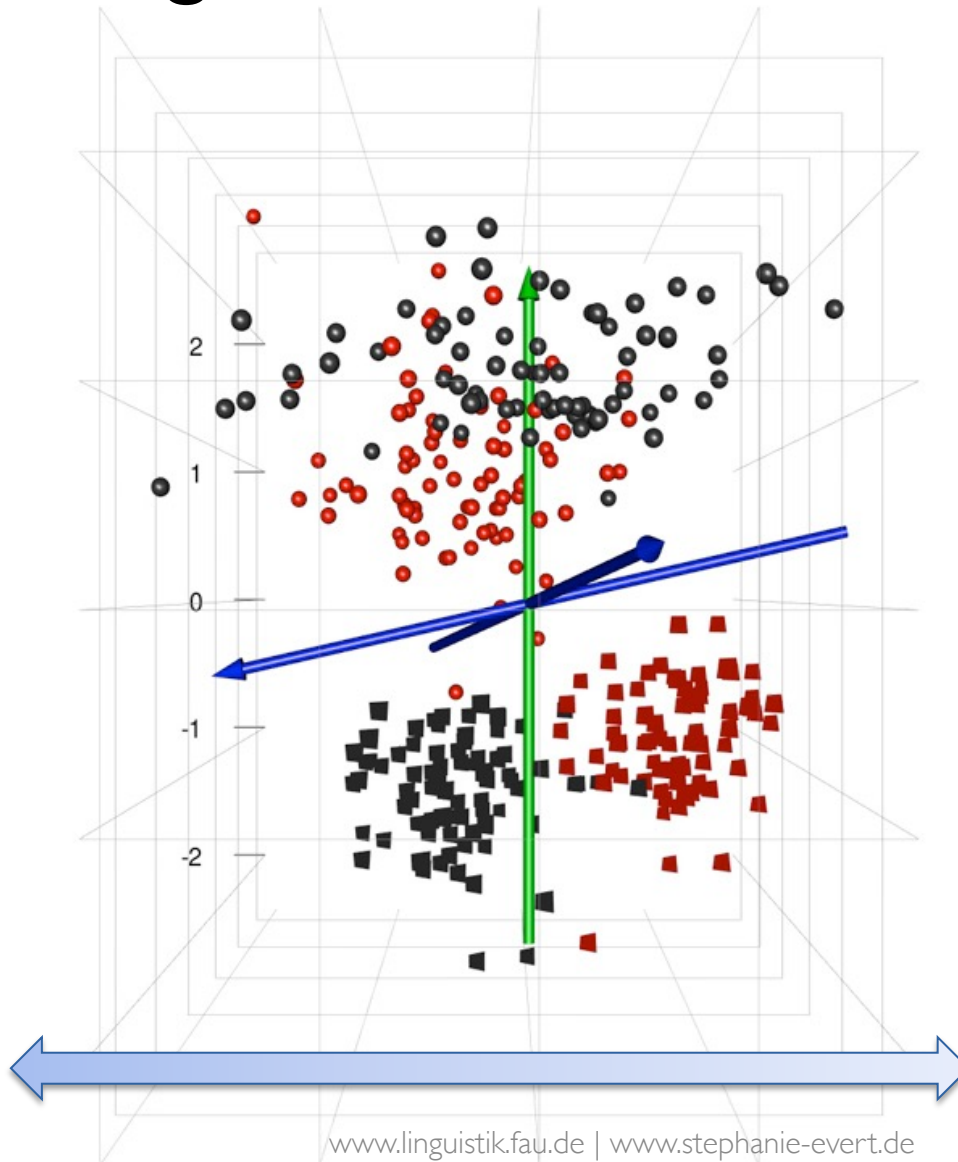
# Unravelling translationese

German



English

SIGIL Unit #7



- DE
- ▲ EN
- orig
- trans

LDA for trans vs. orig  
in each language

# Case study 2: French regional varieties

(Diwersy, Evert & Neumann 2014)

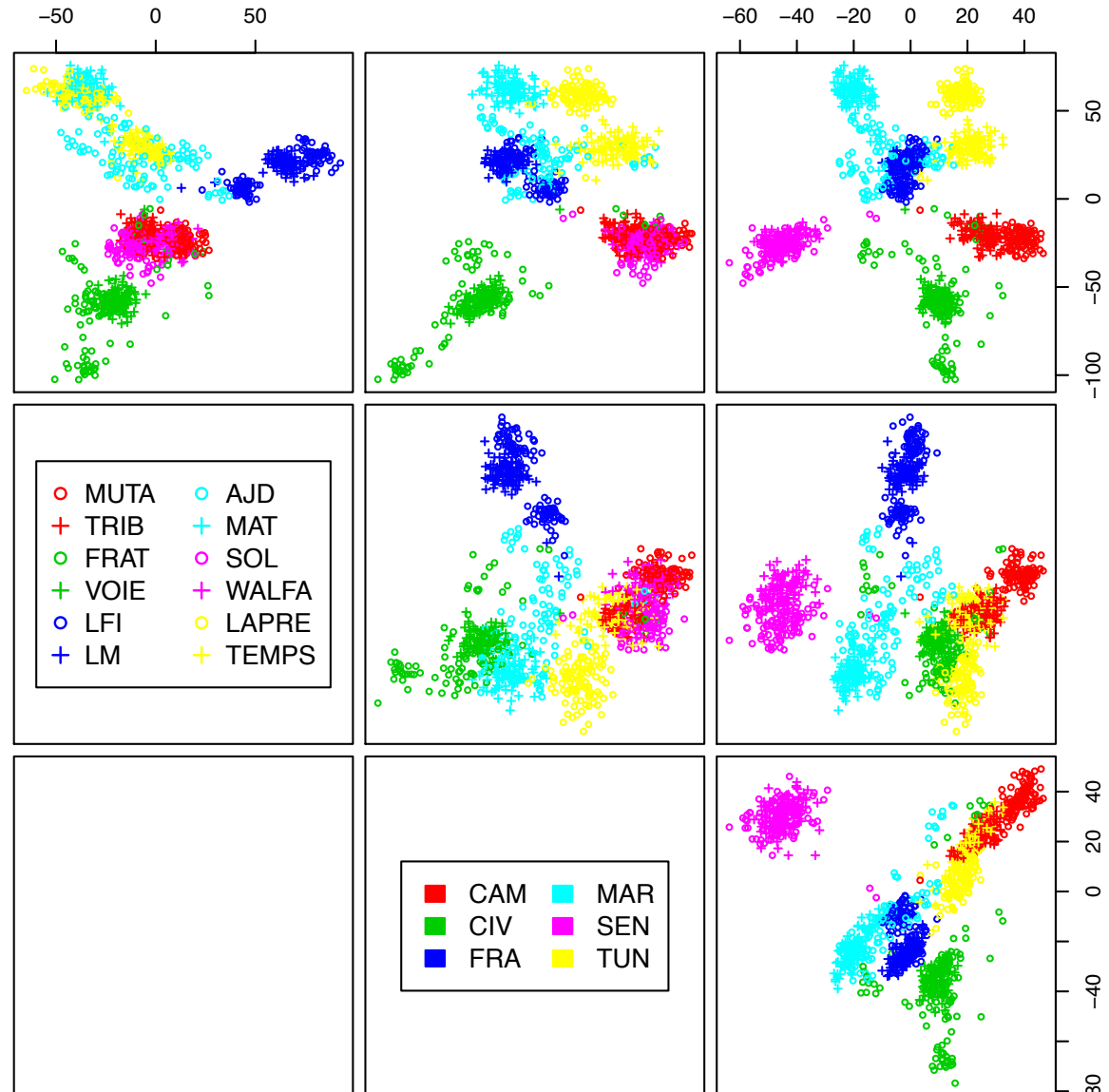
- Lexical differences in regional varieties of French
- Two nation-wide newspapers each from 6 countries
  - Cameroon, France, Ivory Coast, Morocco, Senegal, Tunisia
  - two consecutive volumes from each newspaper
  - total size approx. 14.5 million tokens
- Text samples = one week each
- Features: frequencies of shared colligations
  - colligation = lemma-function pairs
  - must occur in all subcorpora with  $f \geq 100$

# FRV: poor choice of features

PCA **not excluding**  
country-specific  
words as features:  
perfect separation

Design bias results  
in a completely  
uninteresting model

FA not applicable:  
features >> texts

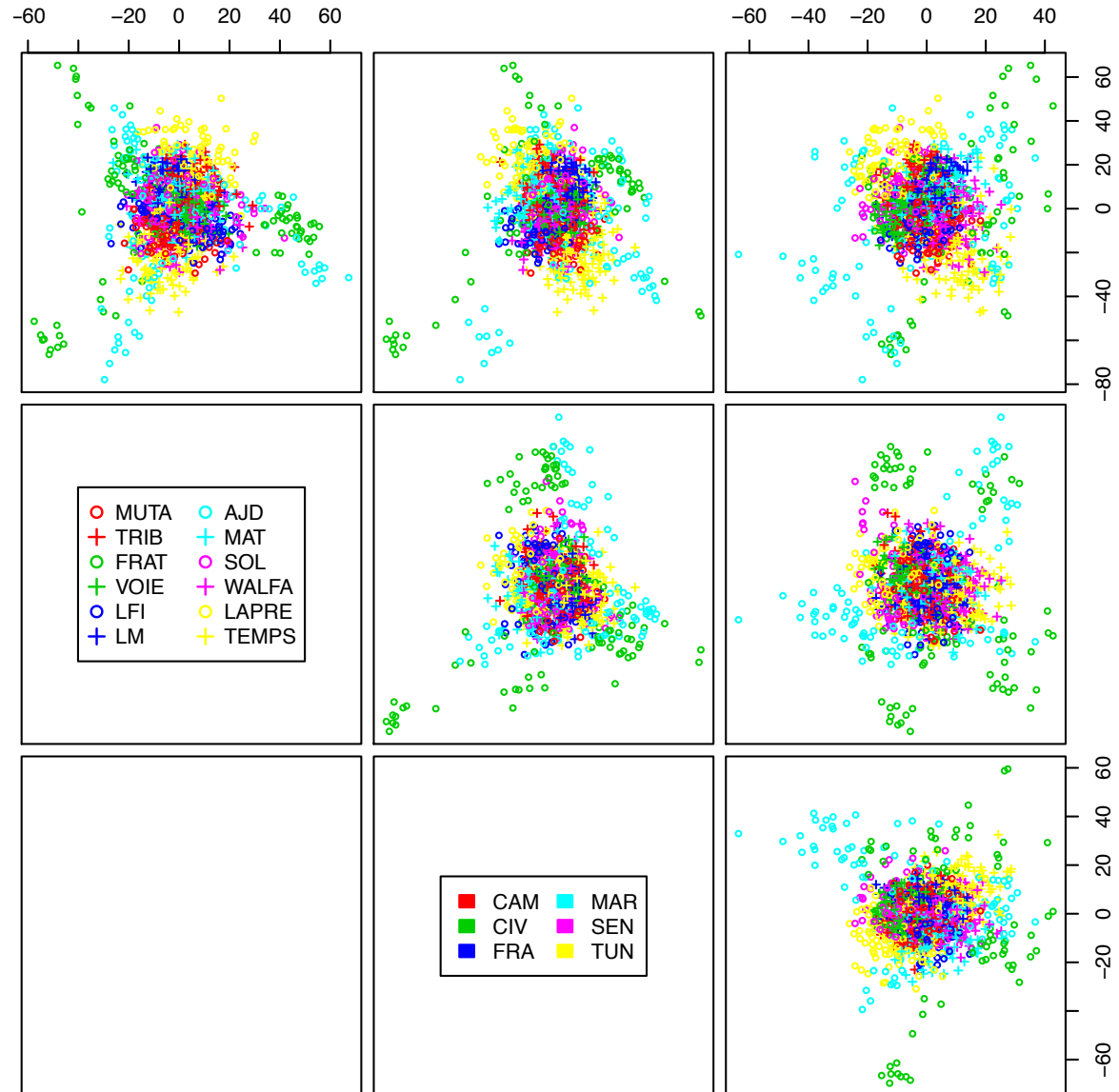


# FRV: PCA dimensions

Using only shared words as features,  
PCA no longer  
reveals any patterns  
(just a few outliers)

Use LDA to find a  
meaningful per-  
spective, based on  
newspaper source

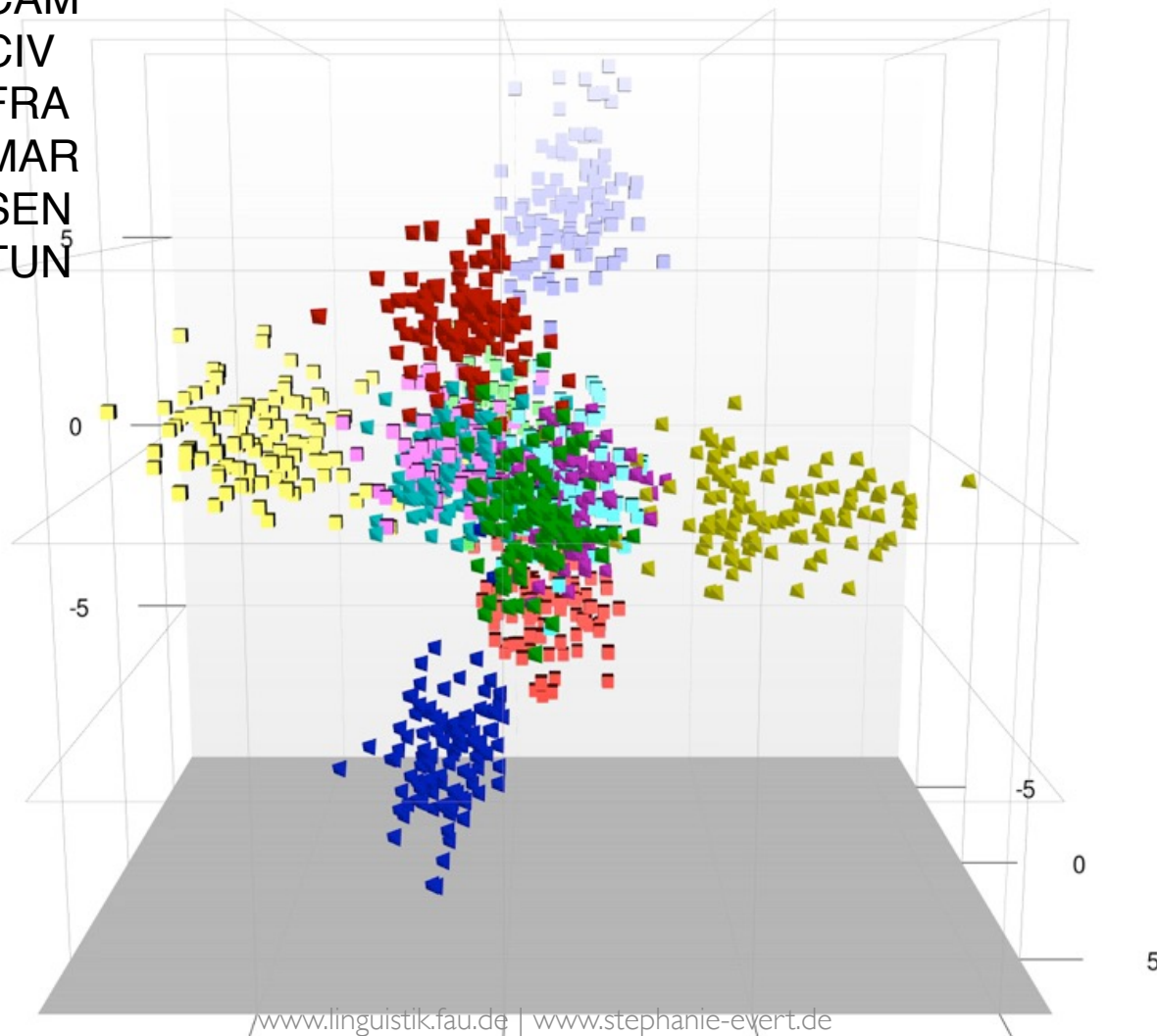
Country would presume  
regional varieties exist!



# FRV: LDA dimensions (newspapers)

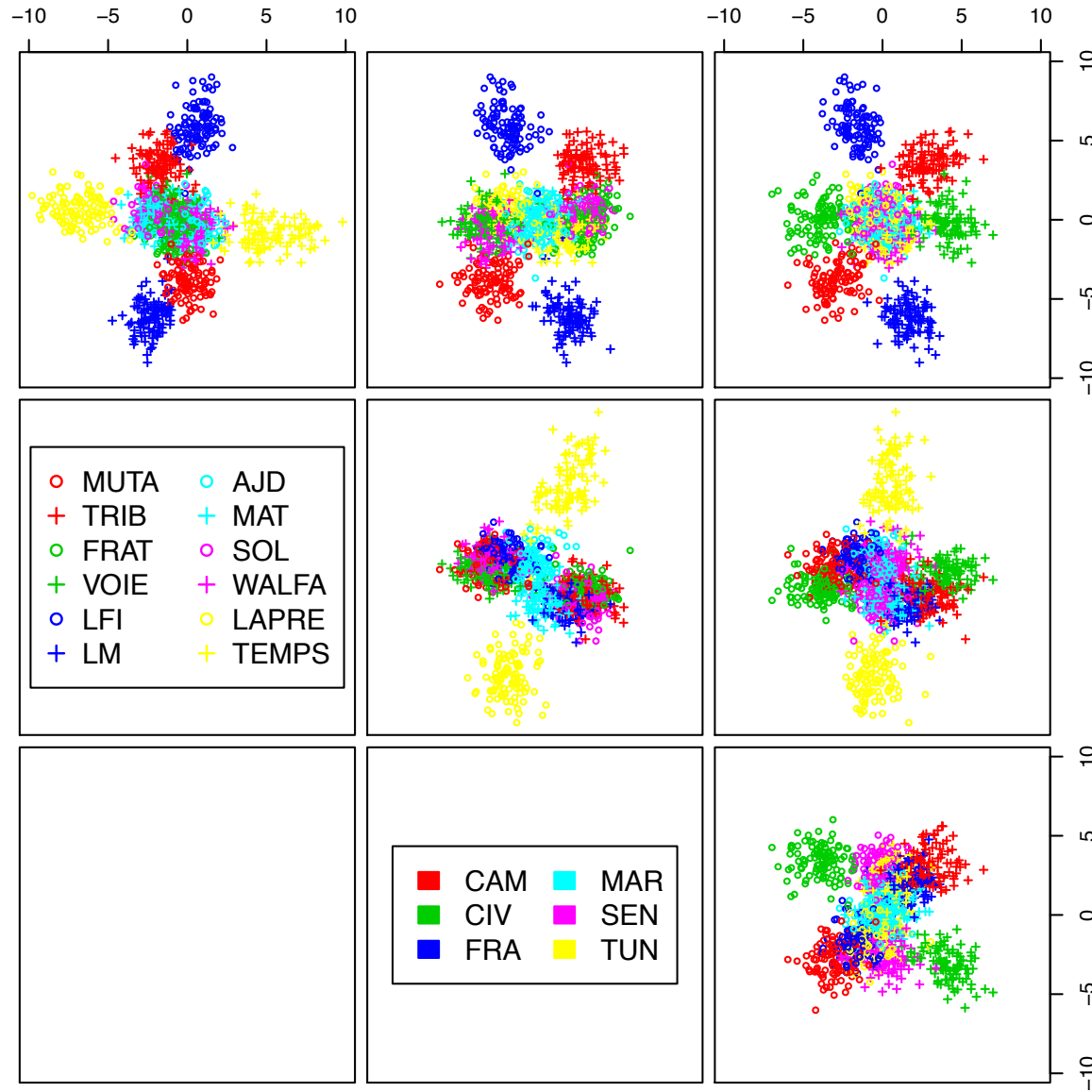
- MUTA
- ▲ TRIB
- FRAT
- ▲ VOIE
- LFI
- ▲ LM
- AJD
- ▲ MAT
- SOL
- ▲ WALFA
- LAPRE
- ▲ TEMPS

- CAM
- CIV
- FRA
- MAR
- SEN
- TUN

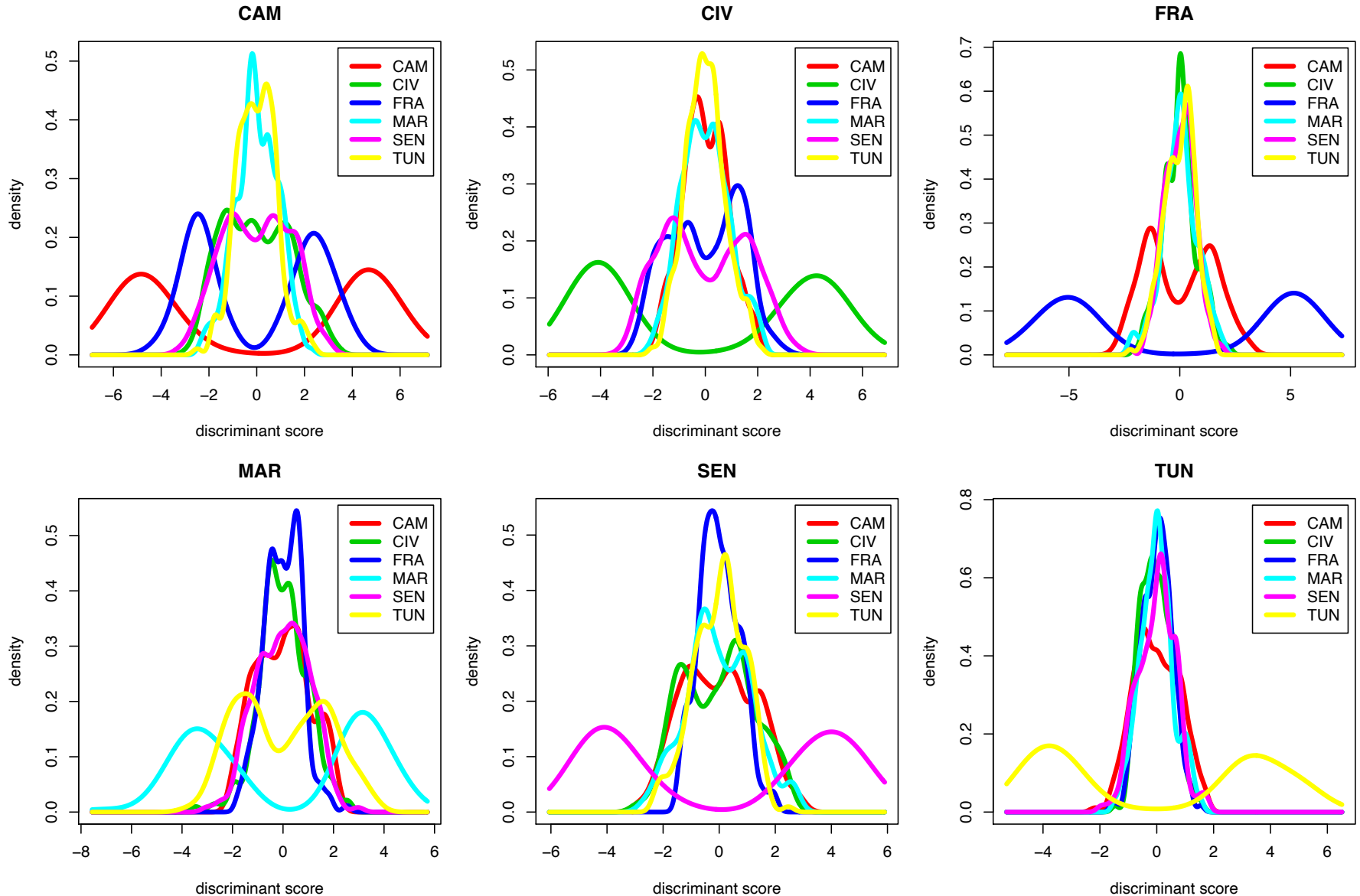




# FRV: LDA dimensions (newspapers)



# FRV: discriminant axes



# References



- Biber, D. (1988). *Variation Across Speech and Writing*. Cambridge University Press, Cambridge.
- Diwersy, S.; Evert, S.; Neumann, S. (2014). *A weakly supervised multivariate approach to the study of language variation*. In B. Szmrecsanyi & B. Wälchli (eds.), *Aggregating Dialectology, Typology, and Register Analysis. Linguistic Variation in Text and Speech*. De Gruyter, Berlin.
- Evert, S. & Neumann, S. (2017). *The impact of translation direction on the characteristics of translated texts: a multivariate analysis for English and German*. In G. De Sutter, M.-A. Lefer & I. Delaere (eds.), *Empirical Translation Studies. New Theoretical and Methodological Traditions (TiLSM 300)*, pages 47–80. Mouton de Gruyter, Berlin.
- Gasthaus, J. (2007). *Prototype-Based Relevance Learning for Genre Classification*. B.Sc. thesis, Universität Osnabrück, Institute of Cognitive Science.
- Koppel, M.; Argamon, S.; Shimon, A. R. (2003). *Automatically categorizing written texts by author gender*. *Literary and Linguistic Computing*, **17**(4), 401–412.
- Neumann, S. (2013). *Contrastive Register Variation. A Quantitative Approach to the Comparison of English and German*. de Gruyter Mouton, Berlin.
- Neumann, S. & Evert, S. (2021). *A register variation perspective on varieties of English*. In E. Seoane & D. Biber (eds.), *Corpus based approaches to register variation*. Benjamins, Amsterdam.
- Teich, E. (2003). *Cross-linguistic variation in system and text. A methodology for the investigation of translations and comparable texts*. Berlin: Mouton de Gruyter.
- Toury, G. (2012). *Descriptive Translation Studies – and beyond: Revised edition*. 2nd ed. Amsterdam: John Benjamins.