



FRIEDRICH-ALEXANDER
UNIVERSITÄT
ERLANGEN-NÜRNBERG

PHILOSOPHISCHE FAKULTÄT
UND FACHBEREICH THEOLOGIE

A multivariate approach to linguistic variation and distribution

Stefan Evert

Corpus Linguistics Group
FAU Erlangen-Nürnberg

Linguistic variation

Variation of a quantitative linguistic feature

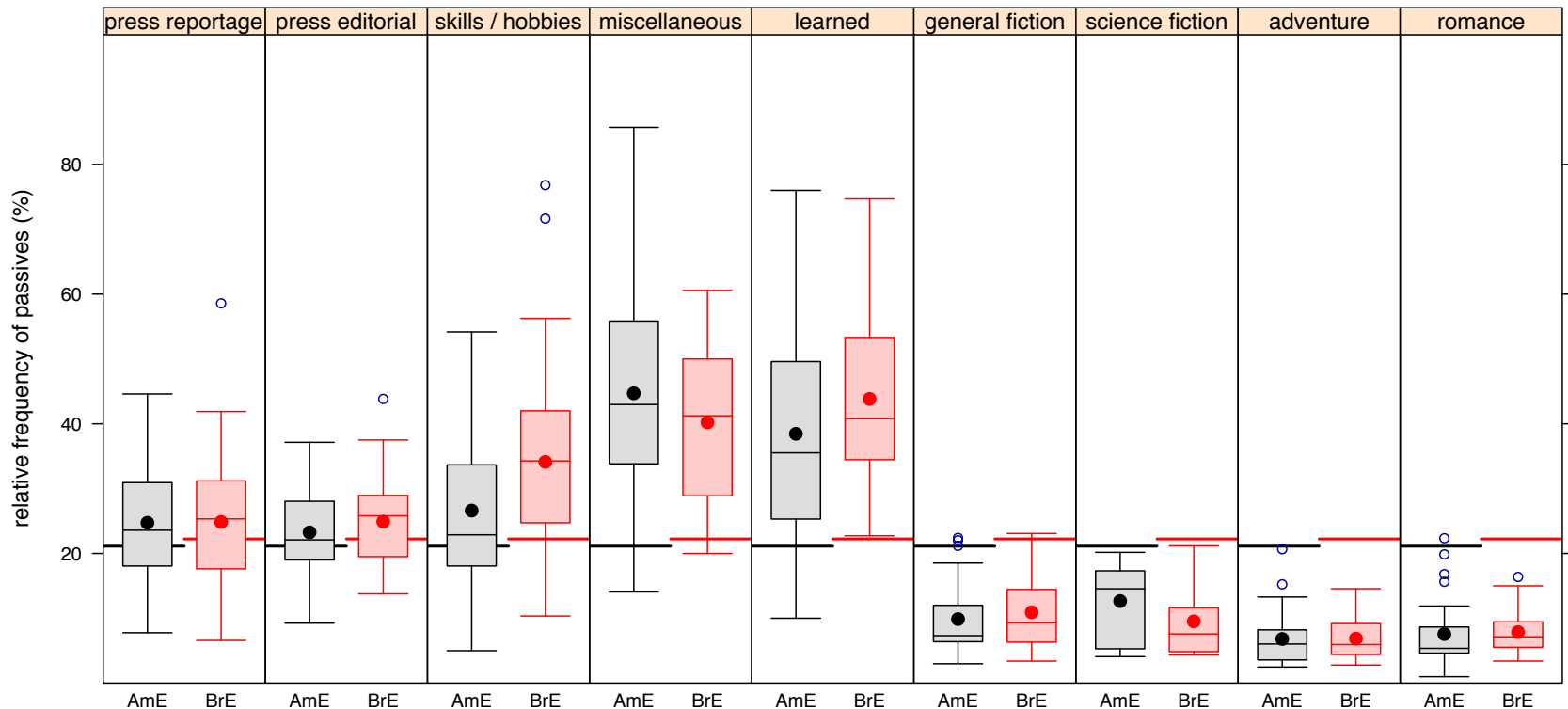
- frequency of passive, past perfect, split infinitive, ...
- frequency of expression, semantic field, topic, ...
- etc.

across

- languages and language varieties
- regions
- social strata
- time
- individual speakers
- etc.

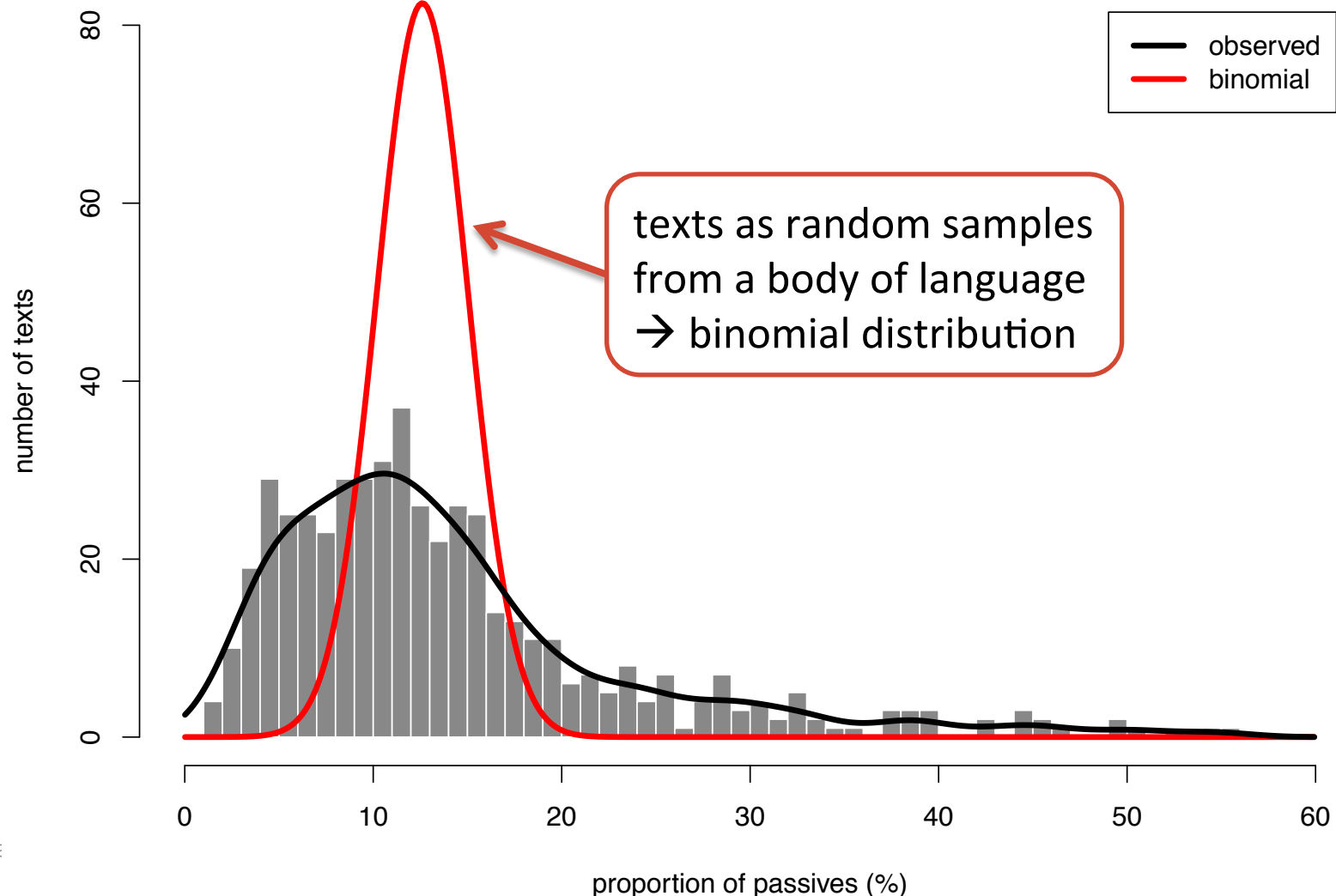
The traditional approach

- Select a linguistic feature (e.g. passive voice)
- Compare its frequency across different categories (genres, language varieties, speakers, ...)



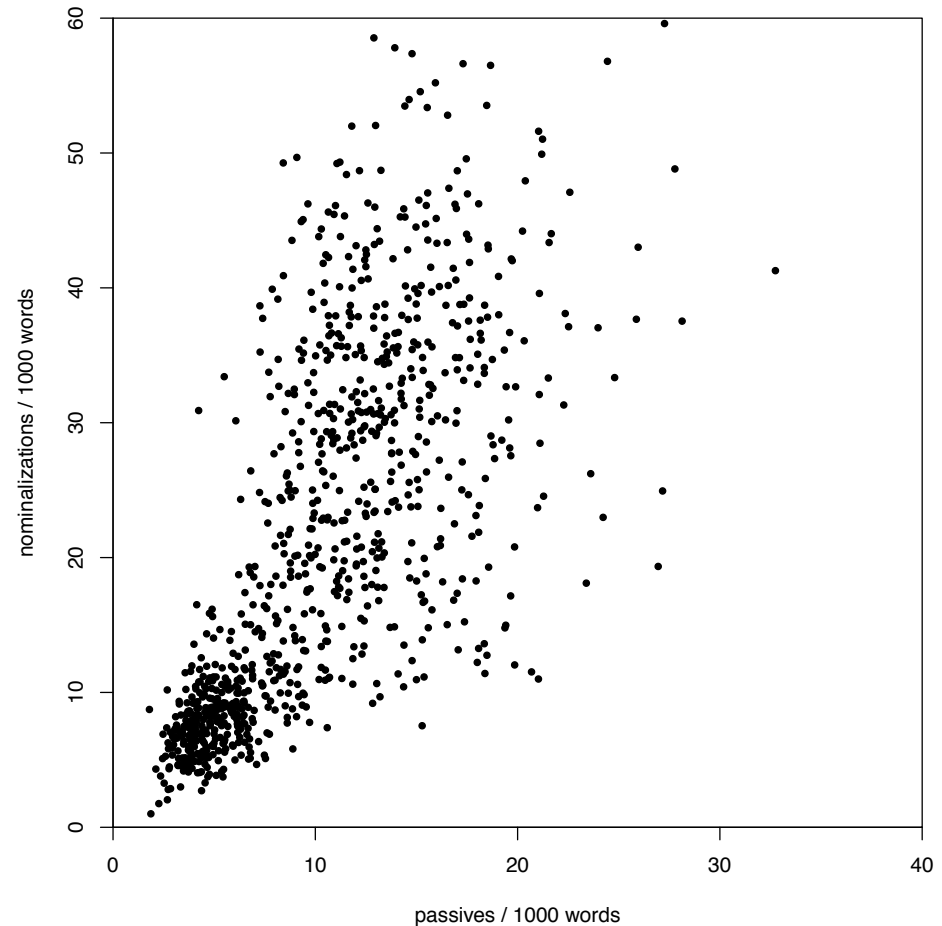
Language variation as a nuisance parameter in corpus linguistics

Passives in Brown corpus (1960s AmE)



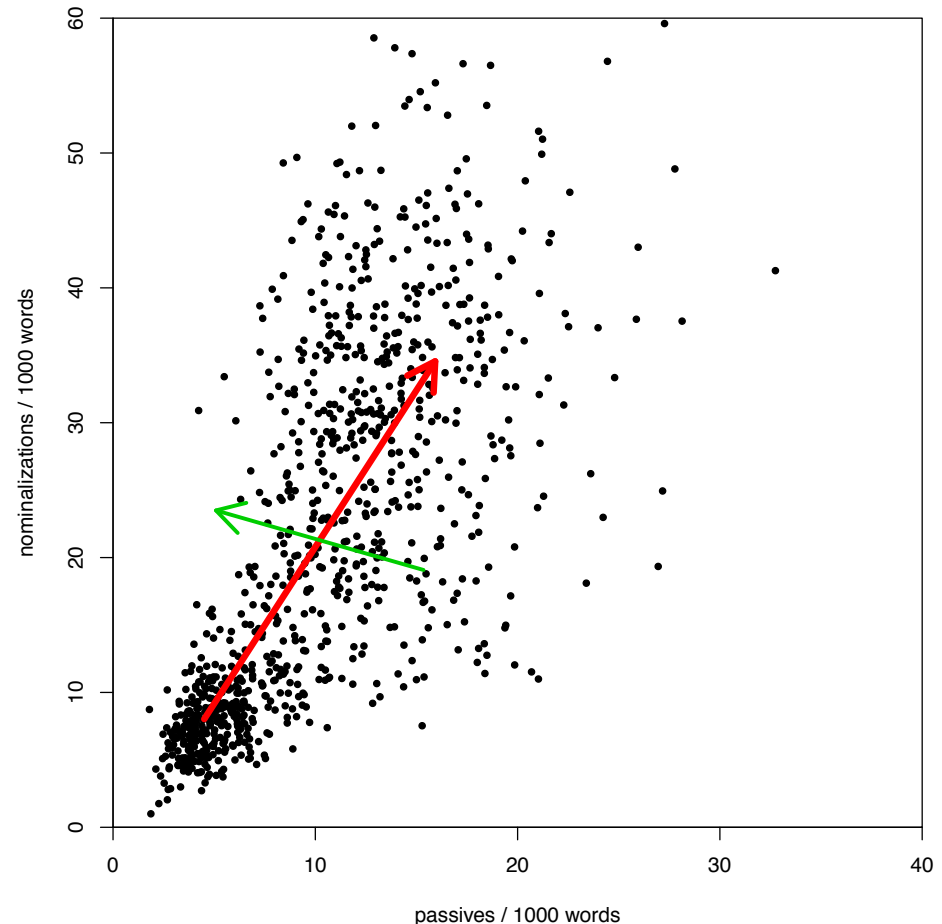
The multivariate approach

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

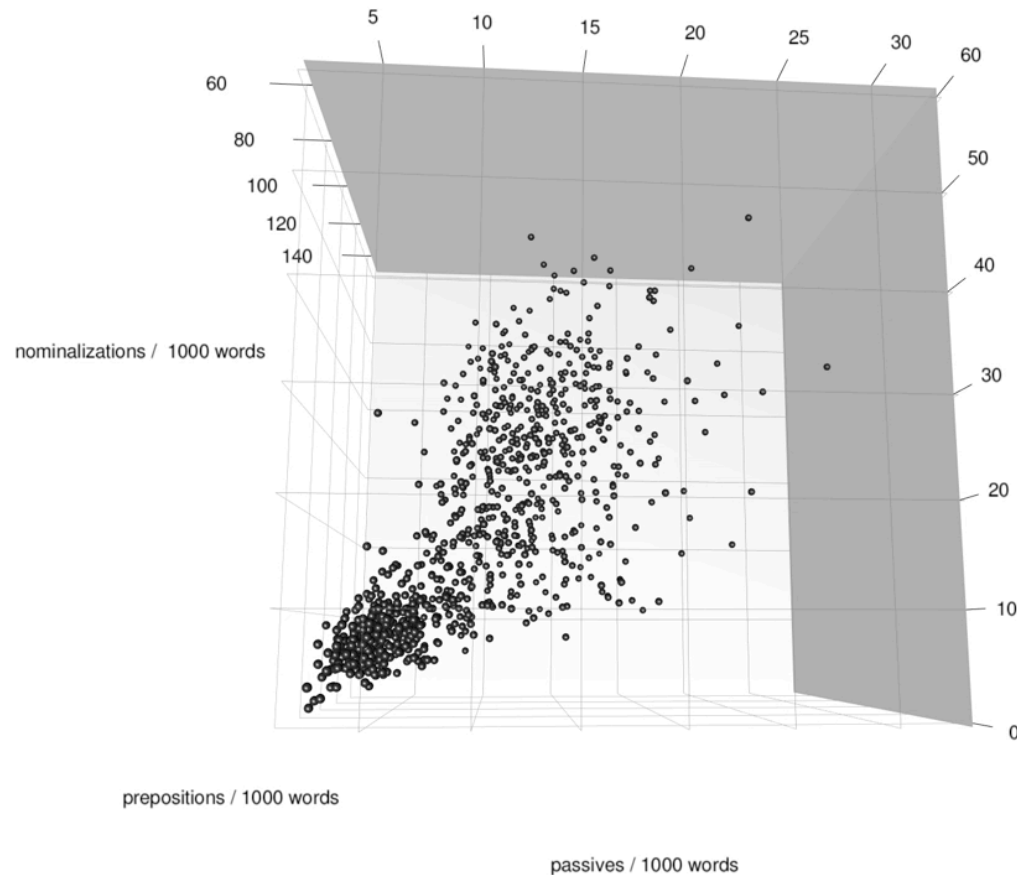


The multivariate approach

- 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

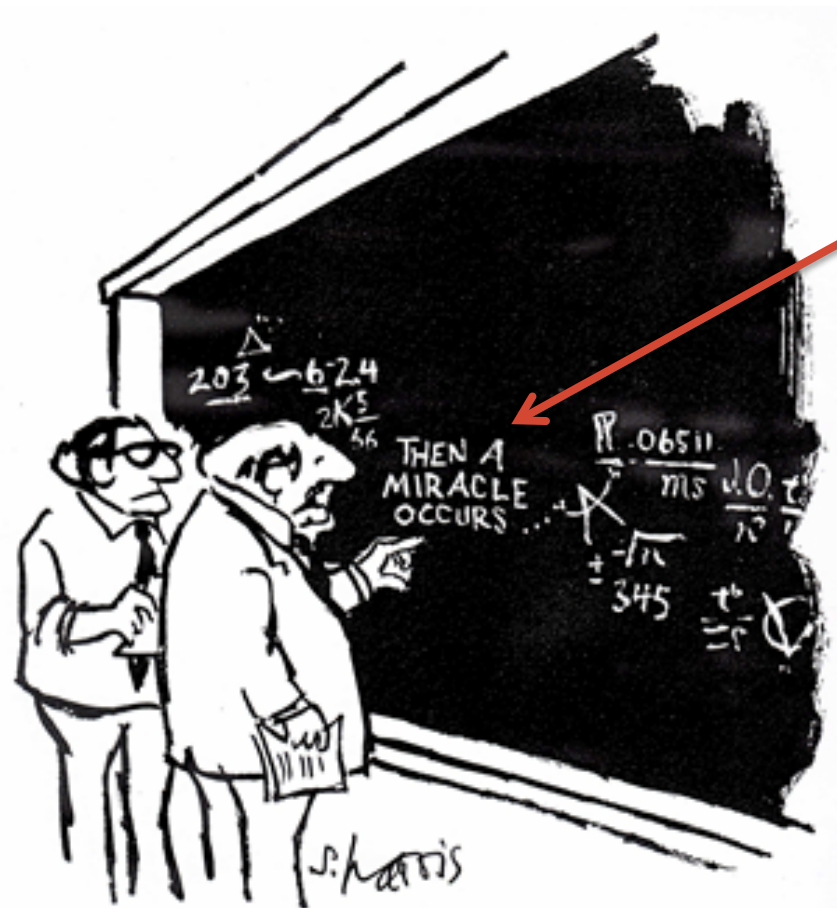
Table 5.7 *Linguistic features used in the analysis of English*

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 multidimensional analysis



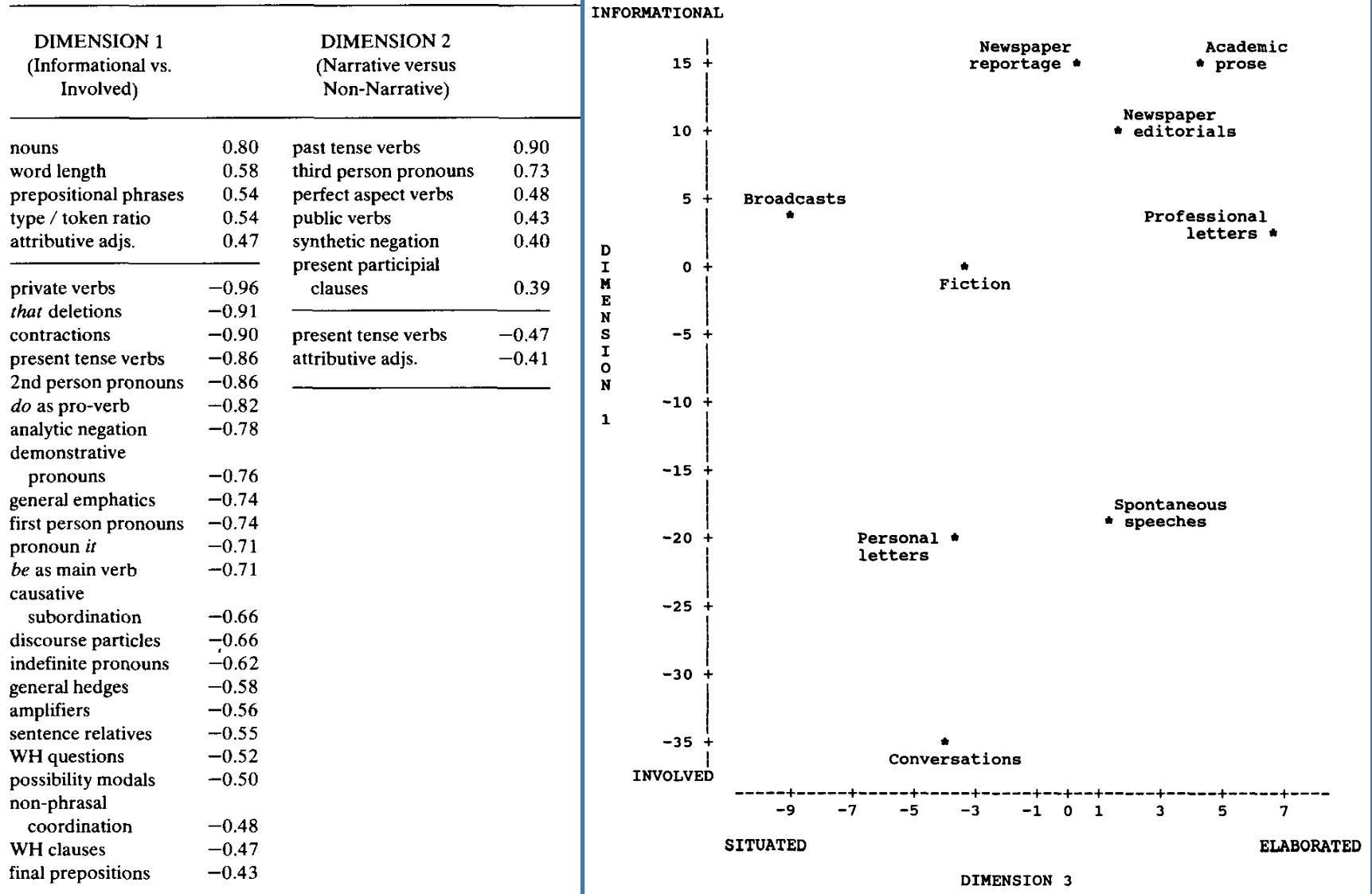
factor analysis
(FA)

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

Biber's multidimensional analysis

TABLE 2

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



Problems

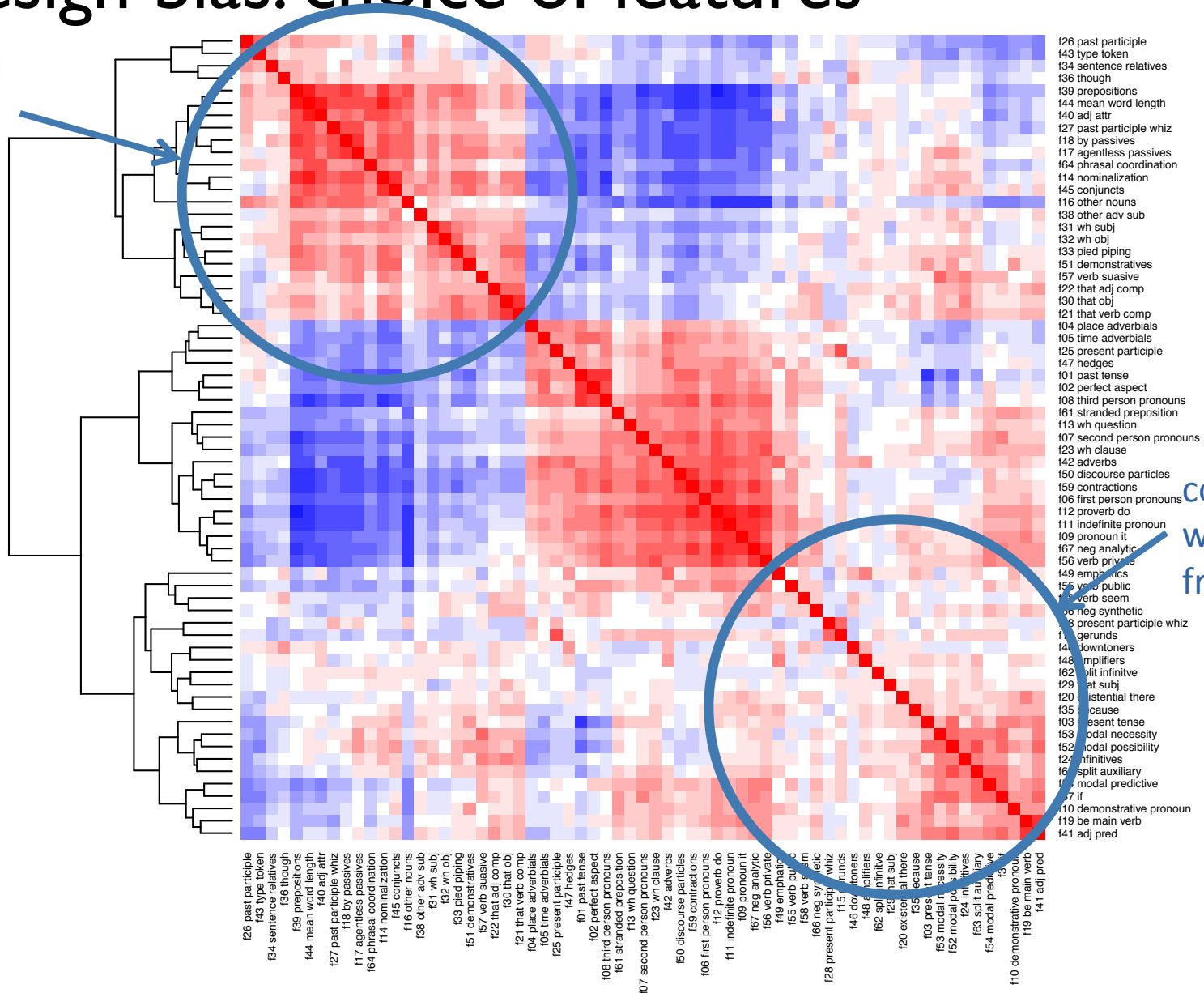
- Design bias
 - choice of features
 - selection of text samples
- Involves a miracle
 - and it isn't even a very robust one
- Interpretation bias
 - arbitrary cutoff for feature weights (“loadings”)
 - risk of reading one's own expectations into features
- More subtle patterns of variation invisible

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)
- Factor analysis with 4 latent dimensions + varimax
 - seems to yield the most clearly structured analysis

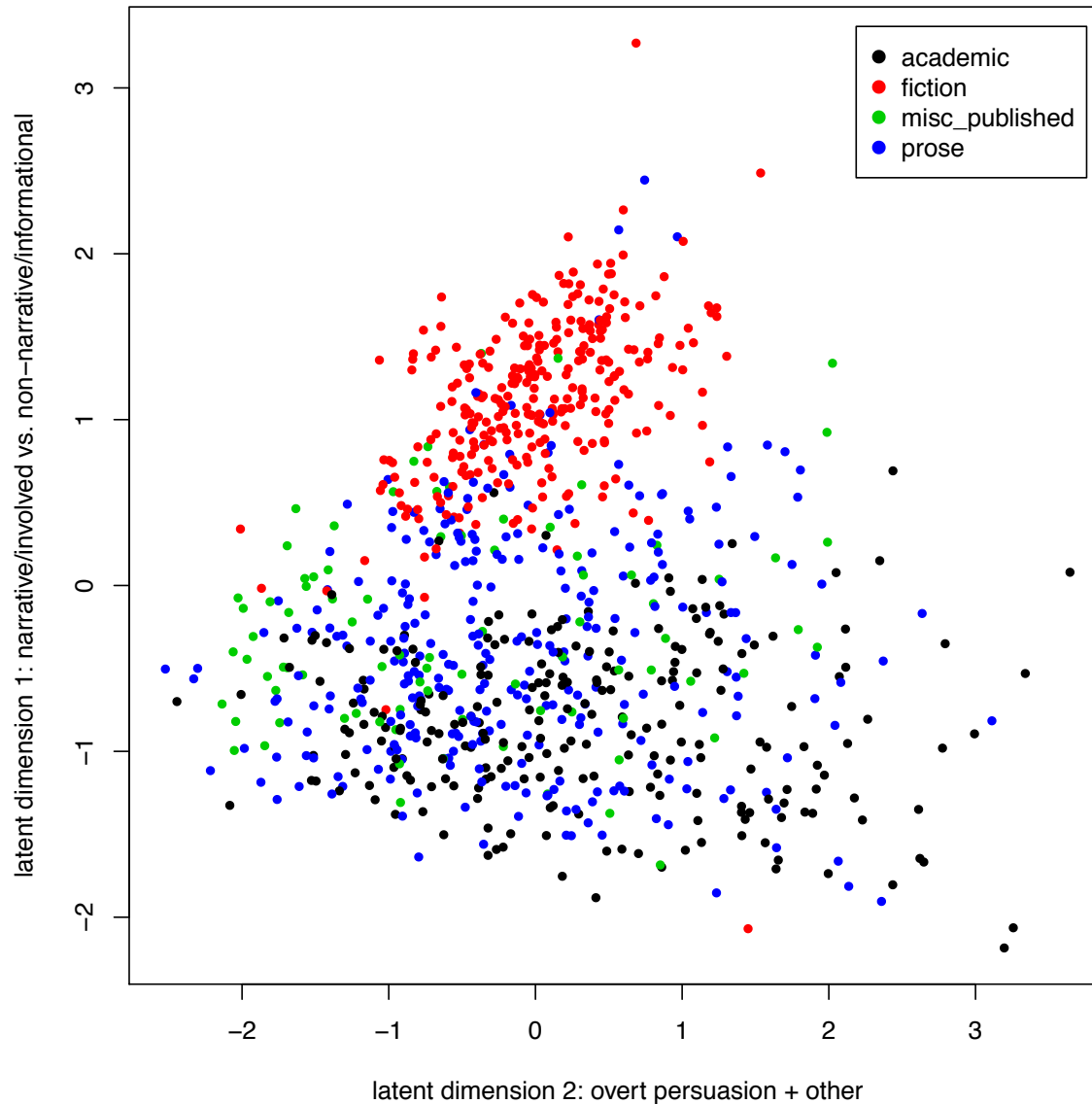
Design bias: choice of features

correlated
with noun
frequency

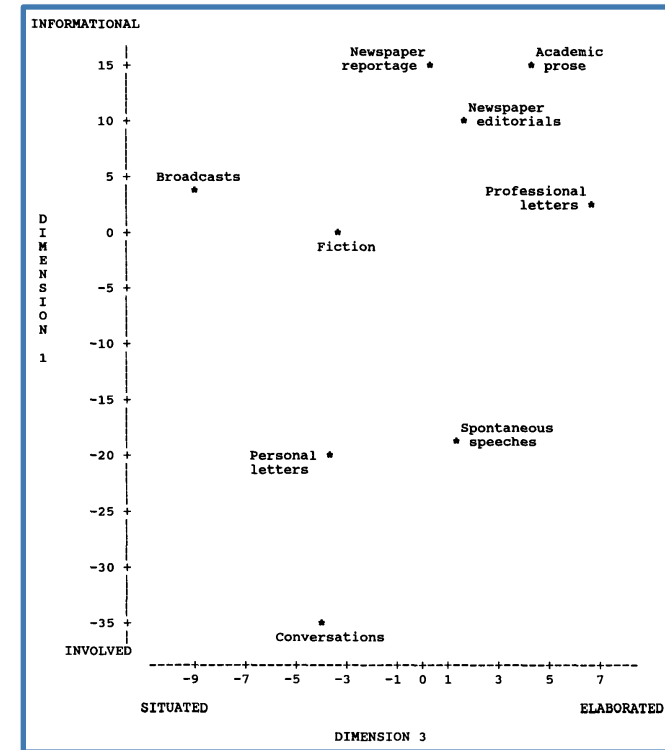
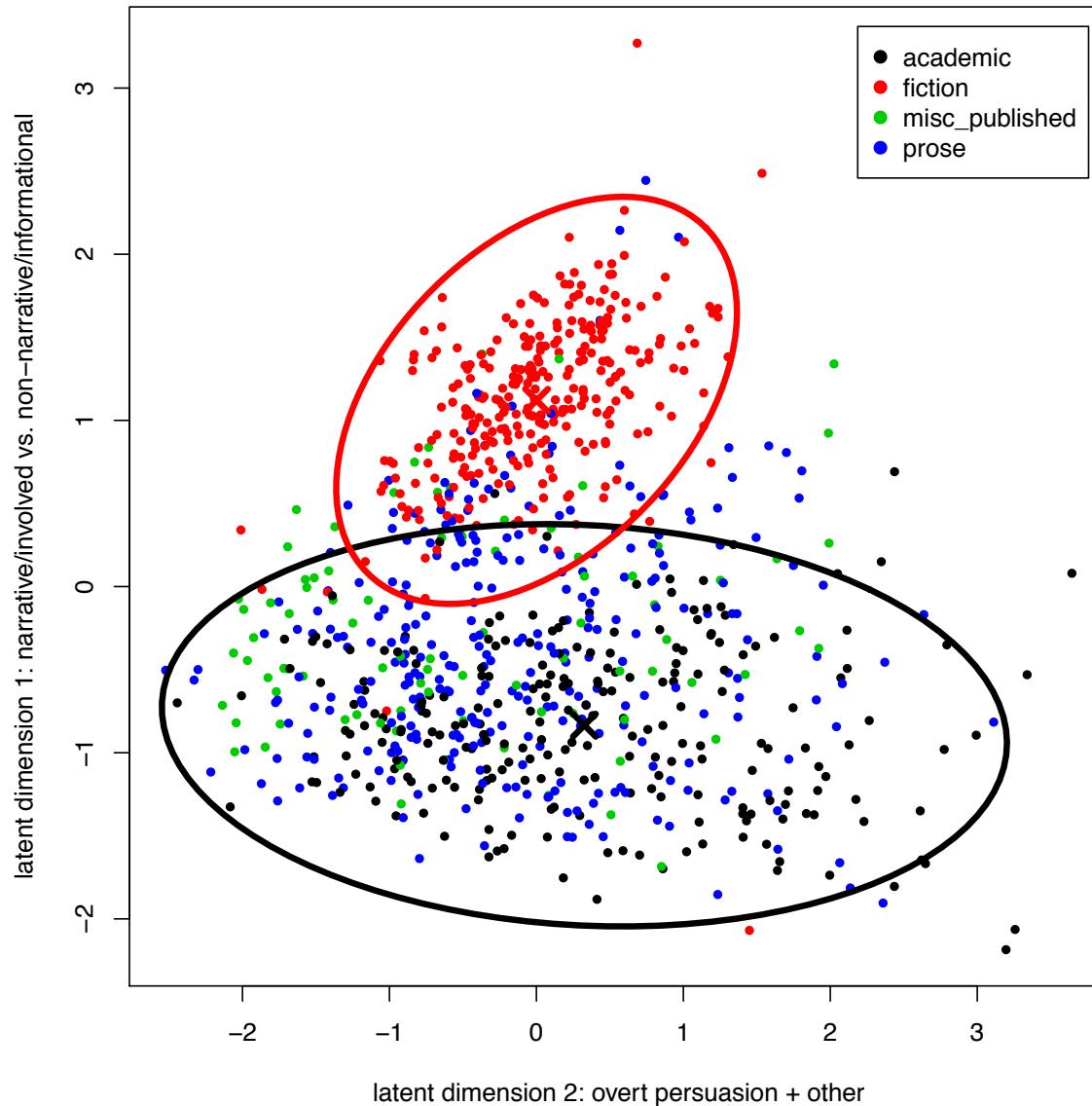


correlated
with verb
frequency

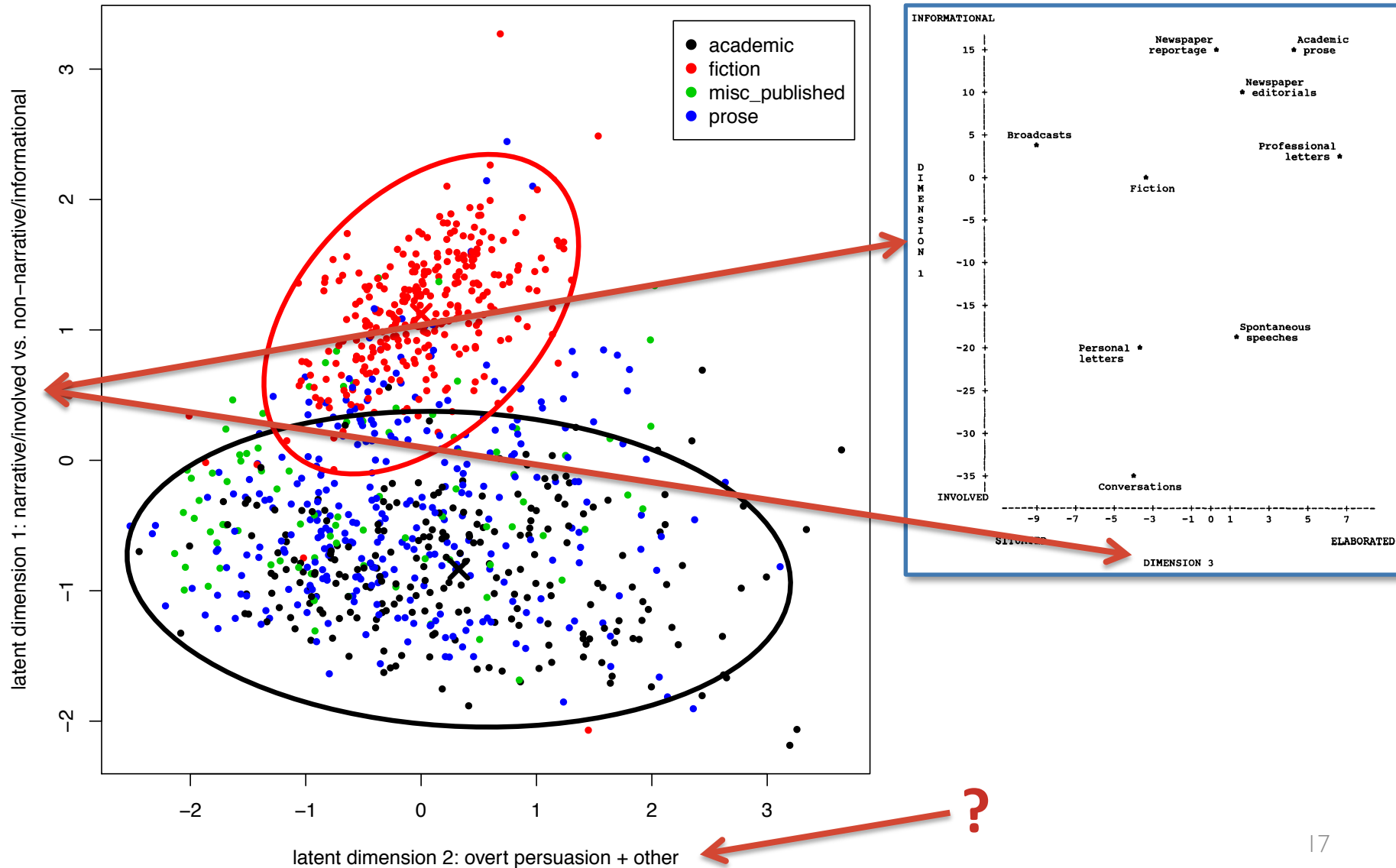
Design bias: choice of texts



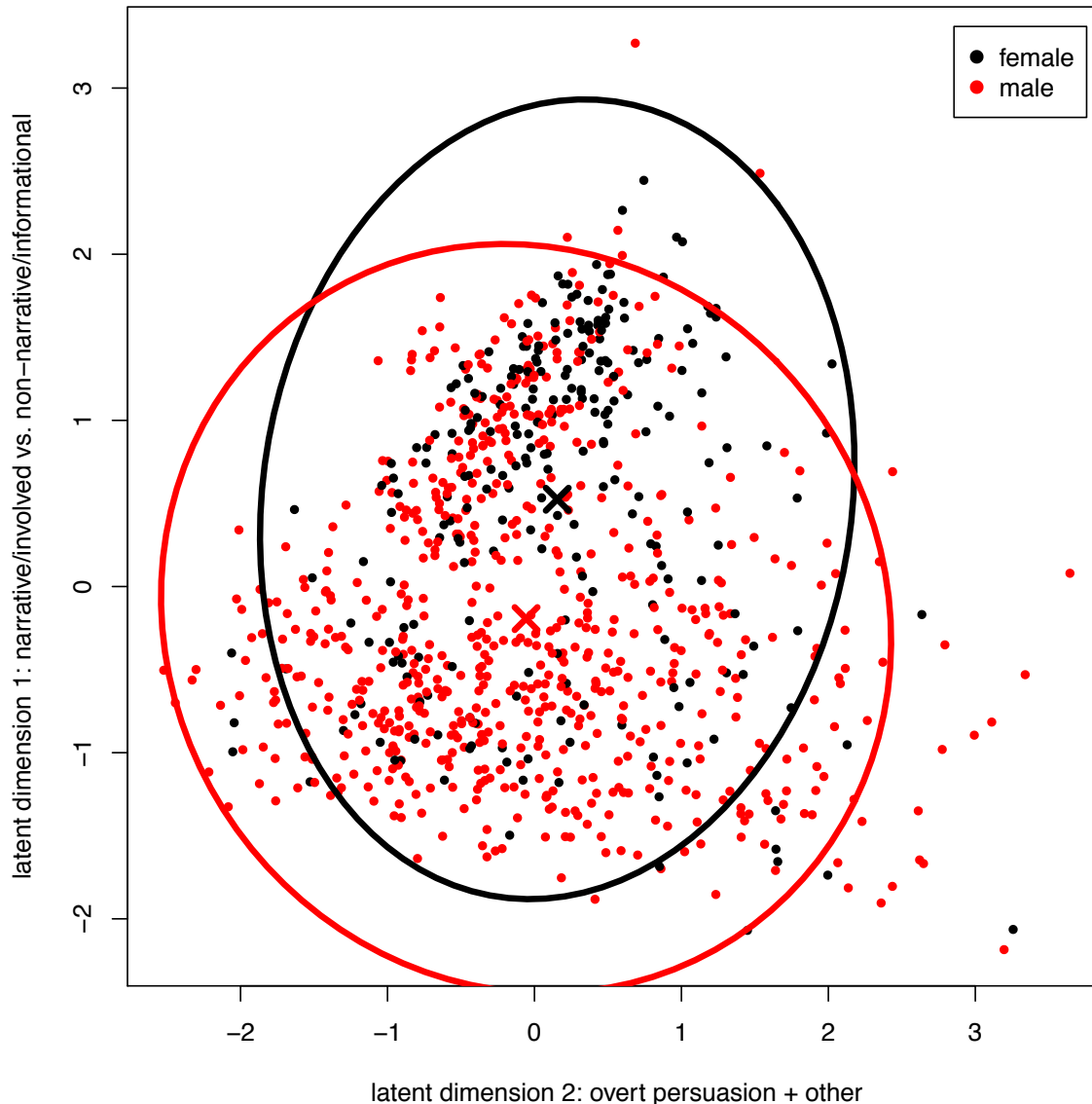
Design bias: choice of texts



Interpretation bias



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%

Our approach

(Diwersy, Evert & Neumann 2014)



FRIEDRICH-ALEXANDER
UNIVERSITÄT
ERLANGEN-NÜRNBERG

PHILOSOPHISCHE FAKULTÄT
UND FACHBEREICH THEOLOGIE



Our approach

(Diwersy, Evert & Neumann 2014)

- Assumption: (Euclidean) distances meaningful
 - as a measure of linguistic similarity of texts
 - depends crucially on choice of features
- Visualization to interpret geometric configuration
- Orthogonal projection = perspective on data
 - (squared) distances decompose into preserved structure + orthogonal (hidden) component
 - optimal projection: principal component analysis (PCA)
- Minimally supervised intervention
 - based on externally observable, theory-neutral information
 - method: linear discriminant analysis (LDA)

Case studies

- Translation effects and register variation in German and English (Evert & Neumann in prep.)
- Regional varieties of French, based on colligational frequencies in newspaper texts (Diwersy et al. 2014)
- Work in progress: Authorship attribution with Burrows Delta (Evert et al. 2015)

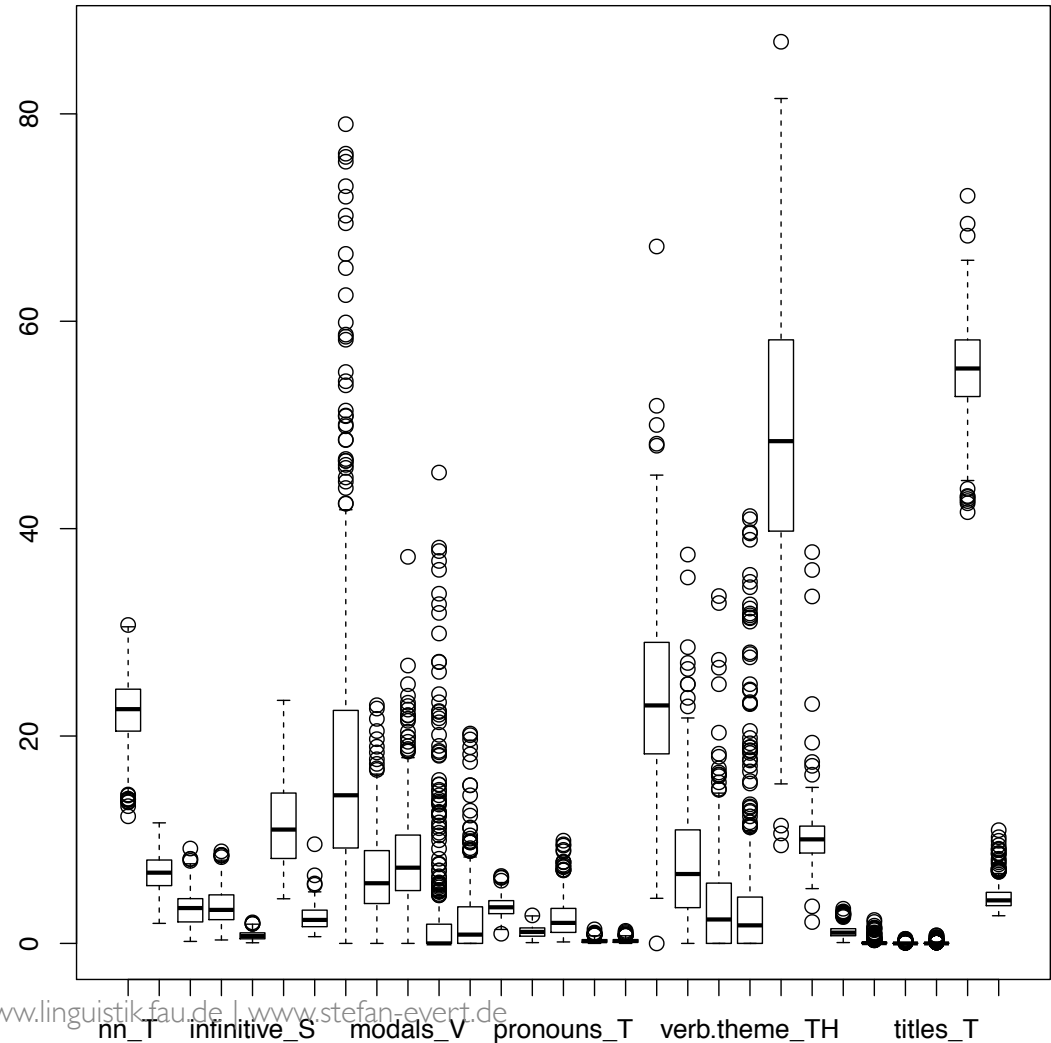
Case study I: CroCo

Diwersy, Evert & Neumann (2014); Evert & Neumann (in prep.)

- CroCo: parallel corpus English/German
 - English-German and German-English translation pairs
 - 454 texts from 8 different genres
- 28 lexico-grammatical features (Neumann 2013)
 - comparable btw. languages, try to reduce correlations
 - inspired by SFL and translation studies
- Text = point in 28-dimensional feature space
- PCA identifies latent dimensions of variation
 - FA results are very similar → comparable to Biber approach
- Focus on English texts here (originals and translations)

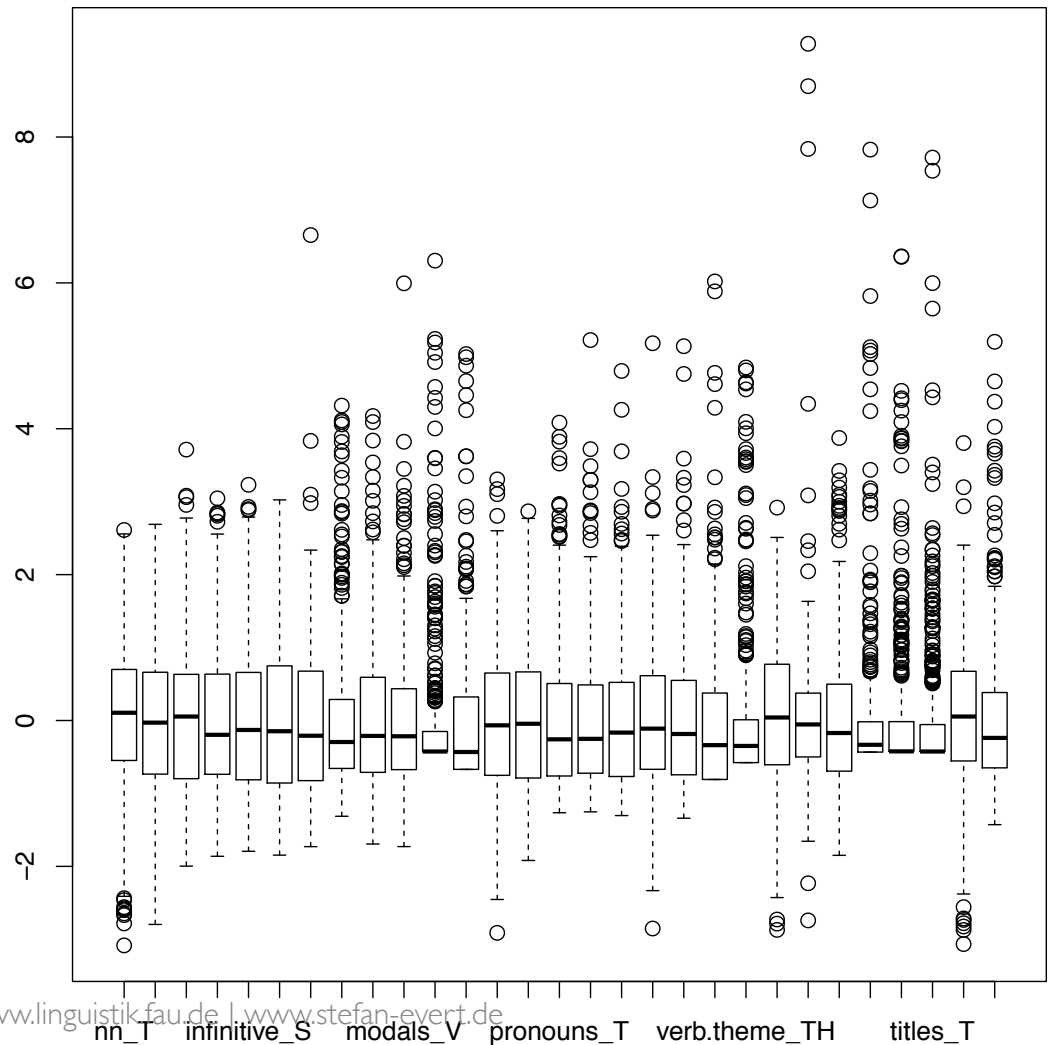
Methodological issues

- Feature scaling
- Choice of features
- Choice of texts
- Delicate effects are obscured

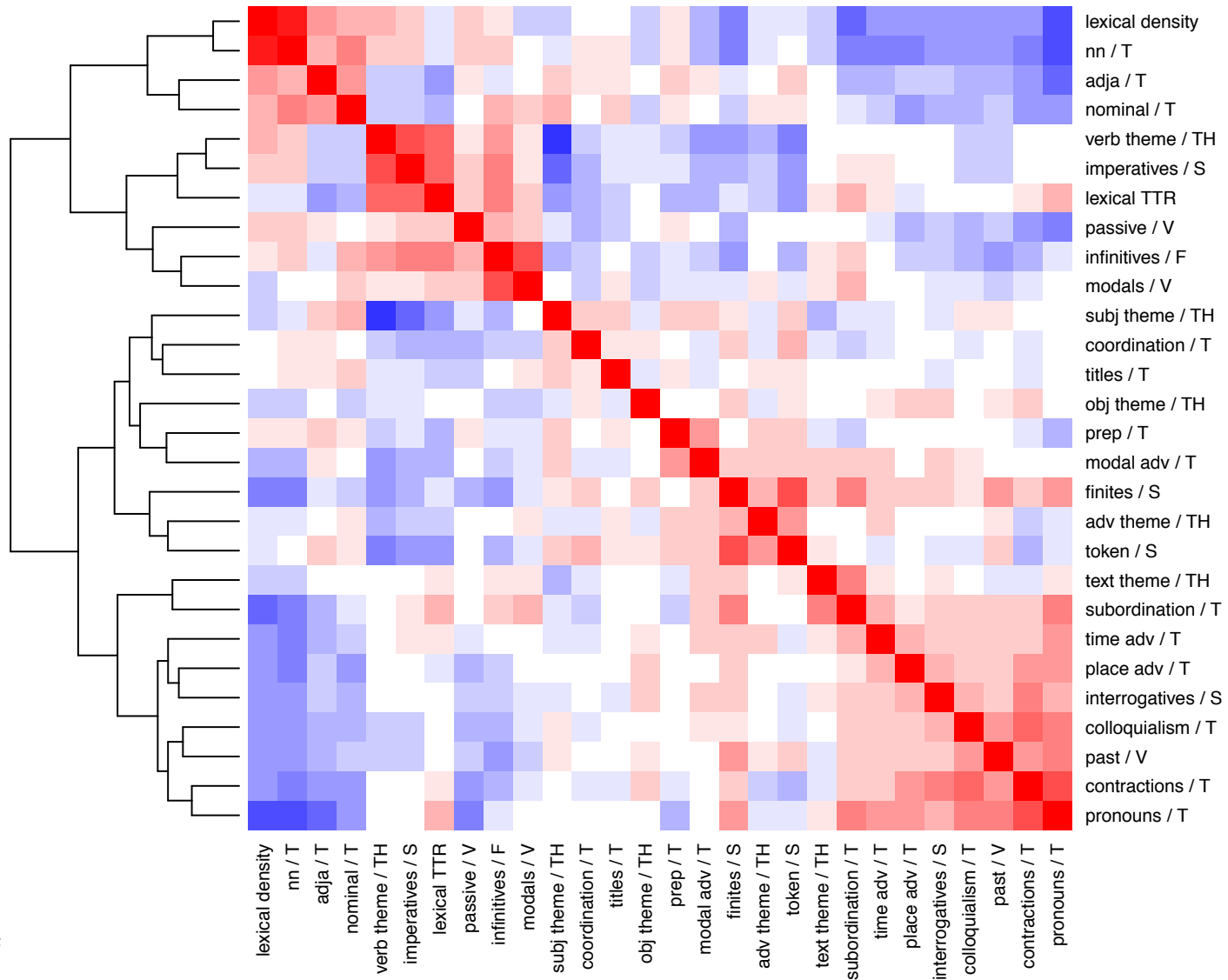


Methodological issues

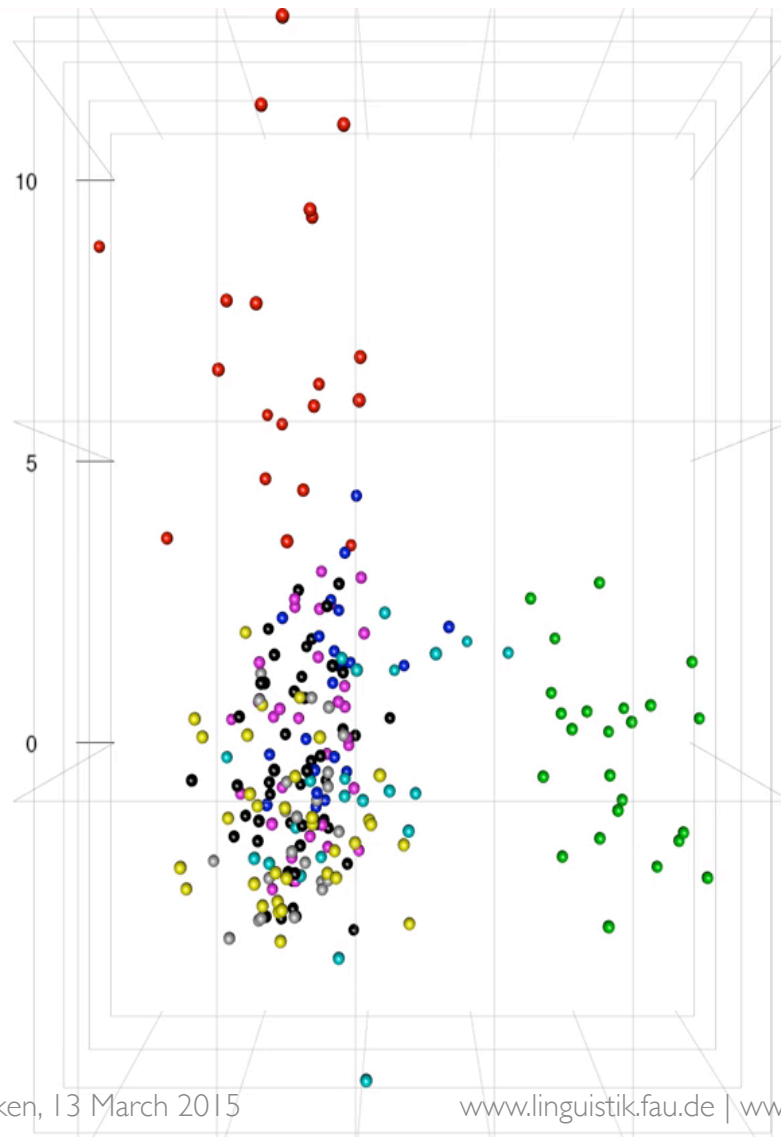
- Feature scaling
- Choice of features
- Choice of texts
- Delicate effects are obscured



Case study I: CroCo

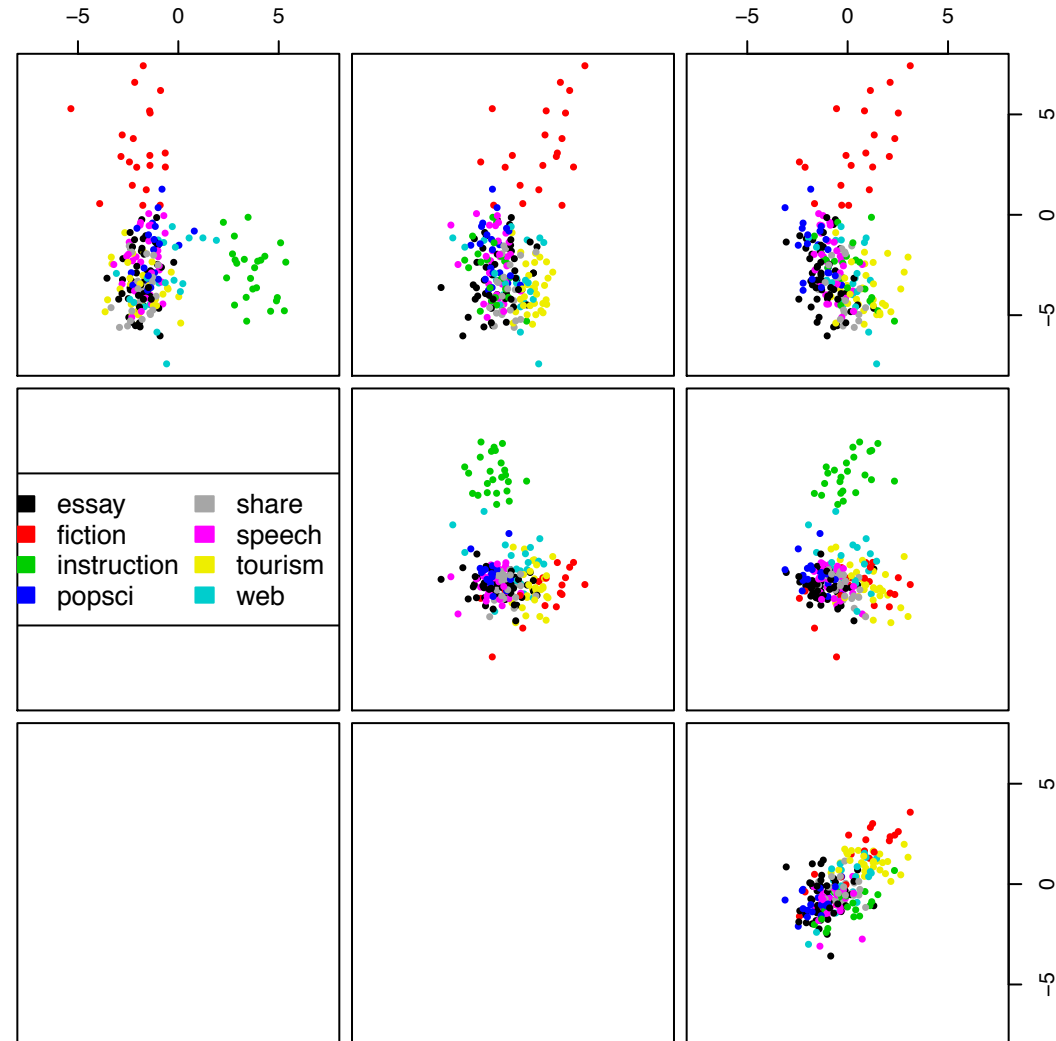
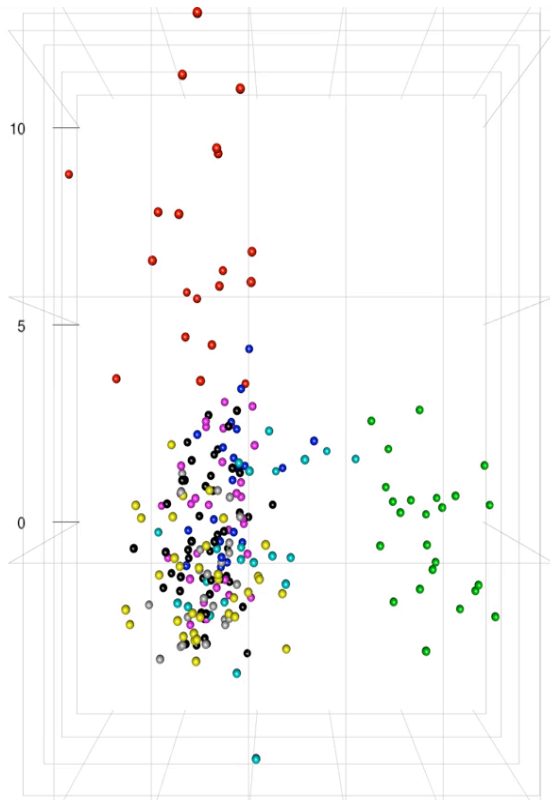


Case study I: CroCo



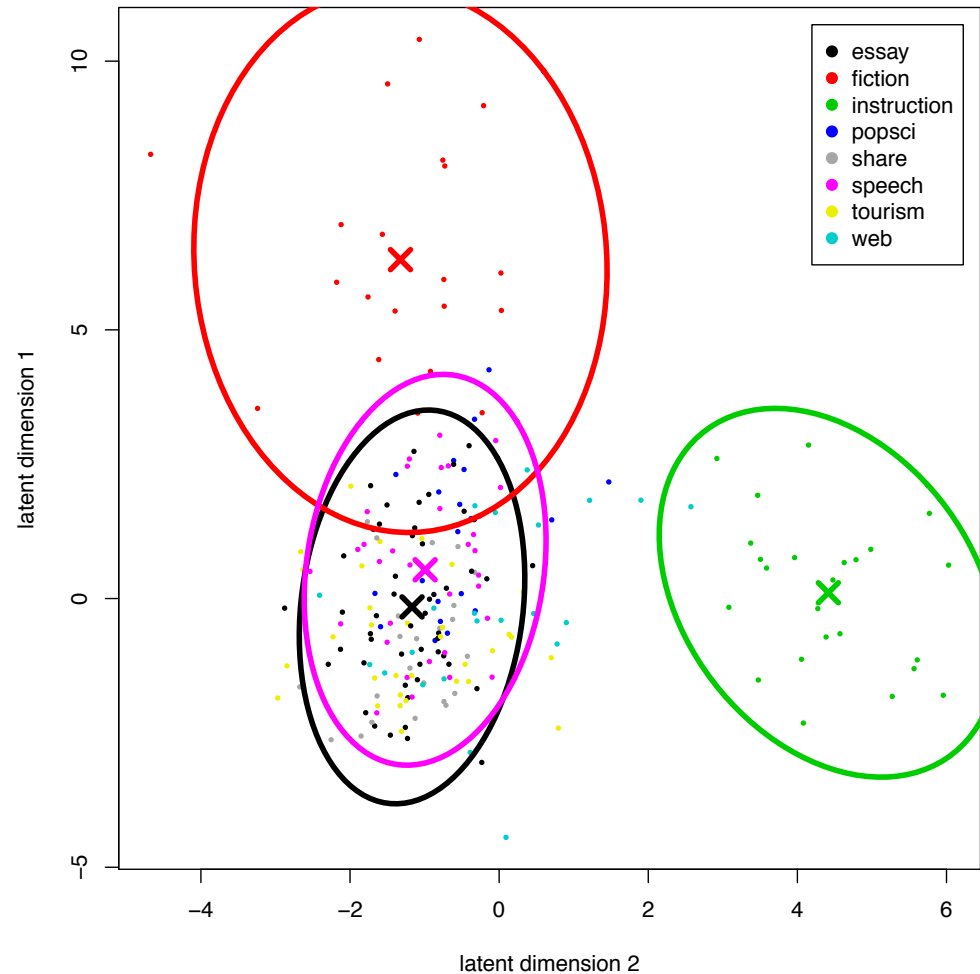
- essay
- fiction
- instruction
- popsci
- share
- speech
- tourism
- web

Case study I: CroCo



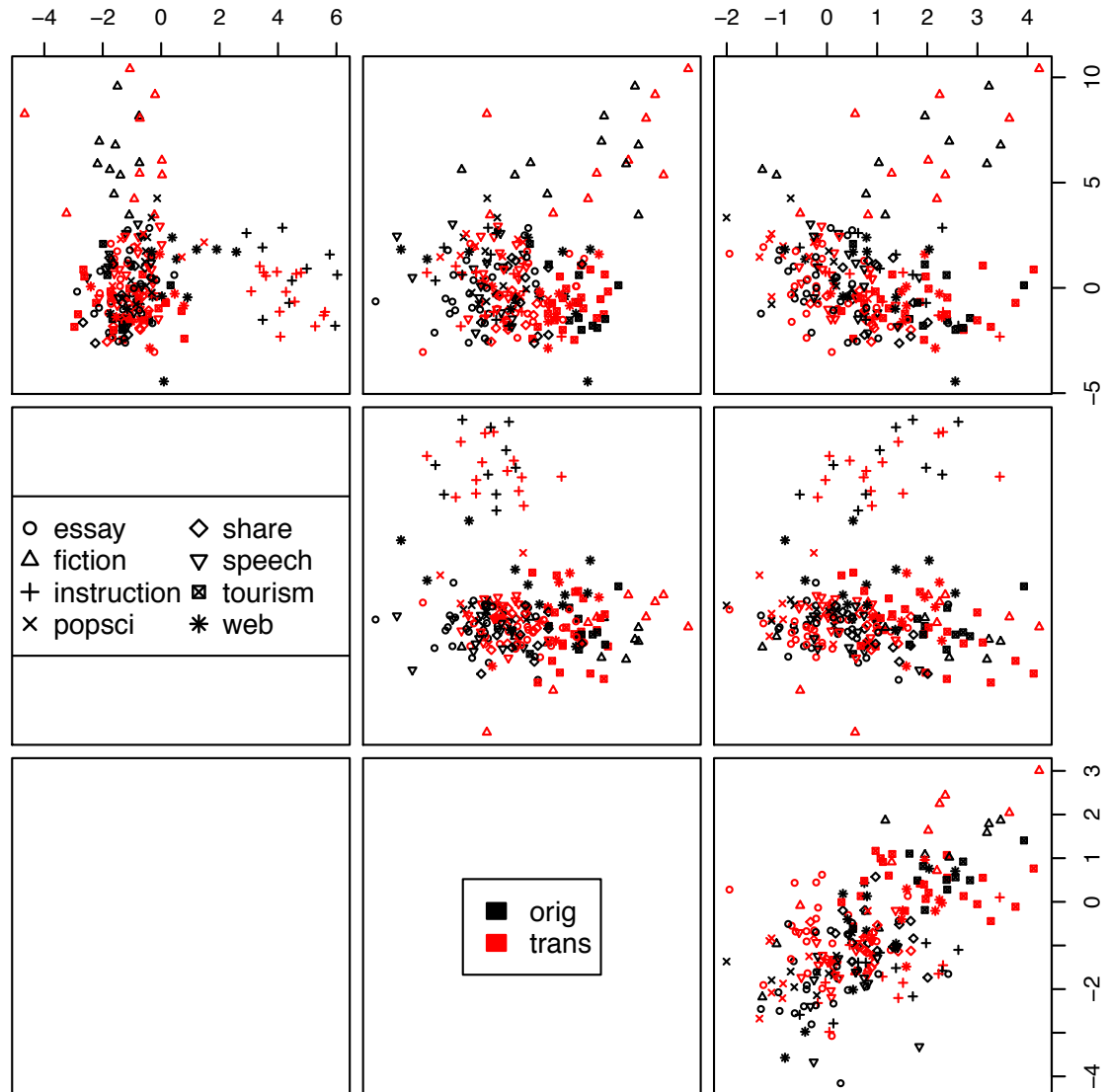
Case study I: CroCo

- Focus on first two latent dimensions (→ Biber's map)
- Describe genre by centroid and confidence ellipse
- Comparison with Hotelling's t^2 test
 - essay vs. speech
 - $t^2=2.512$, $df=2/80$, $p=.0875$ **n.s.**

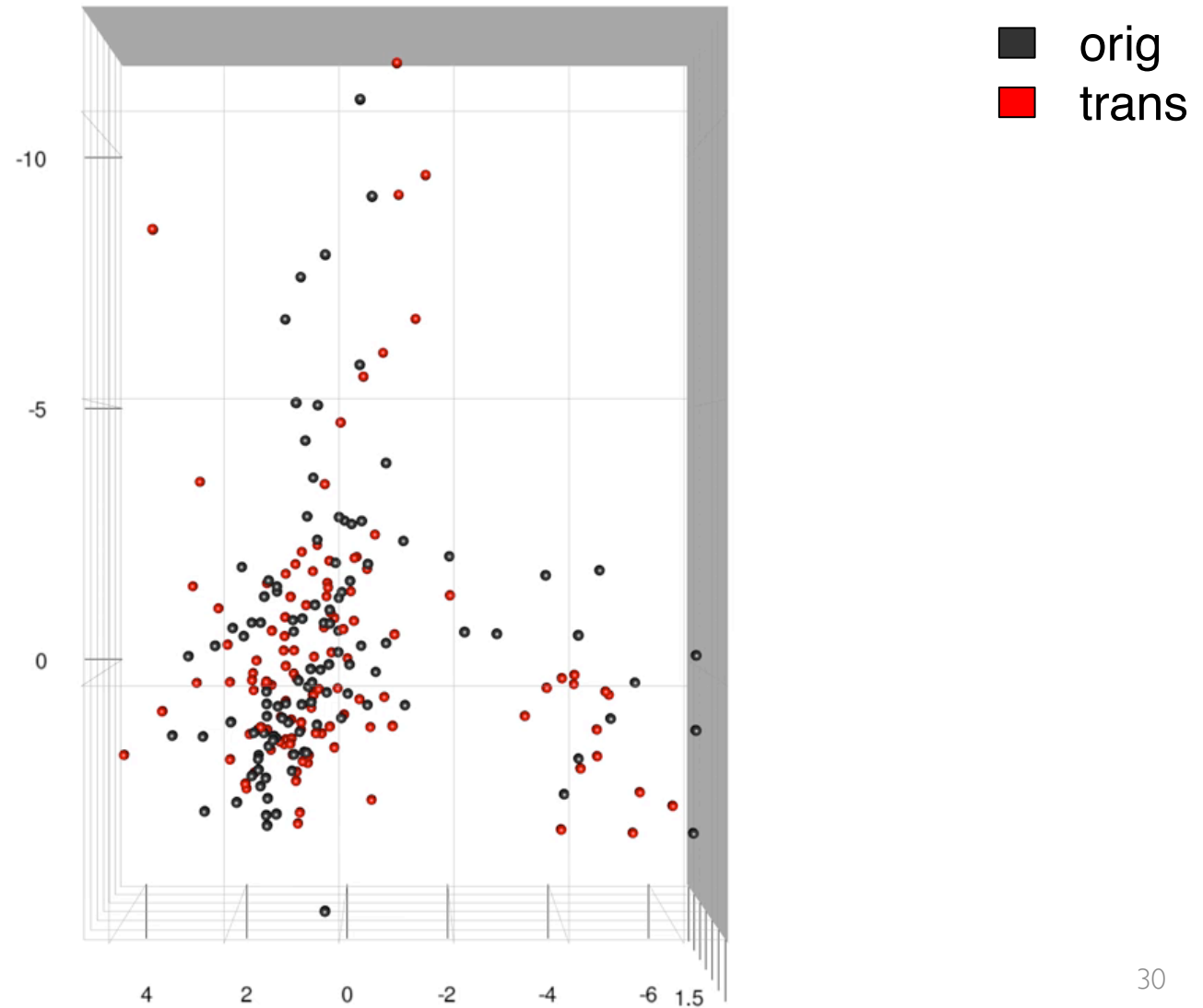


How about subtle patterns?

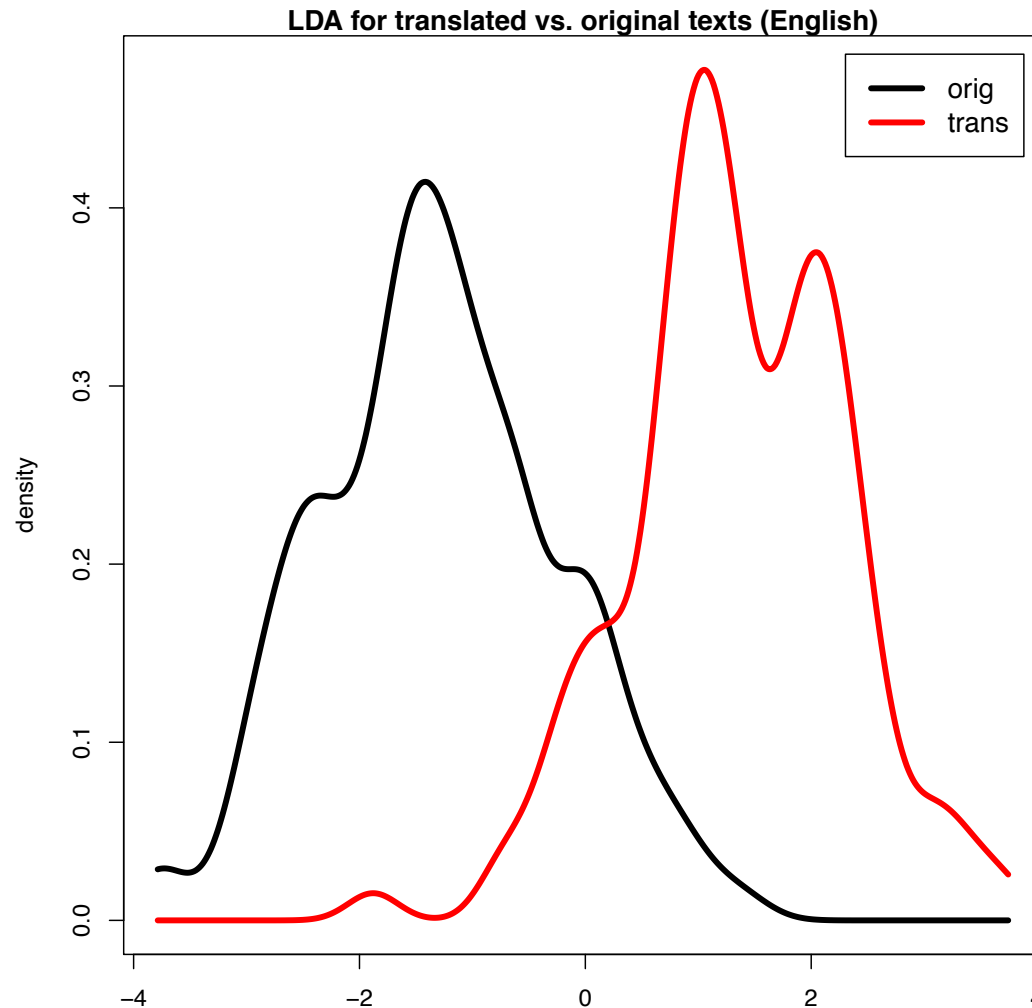
- PCA dimensions fail to distinguish translations from original texts
- But a SVM machine learner can do this with 85% accuracy
- Replace one PCA dimension with LDA discriminant for **orig** vs. **trans**
 - external & theory-neutral information



Finding the right perspective



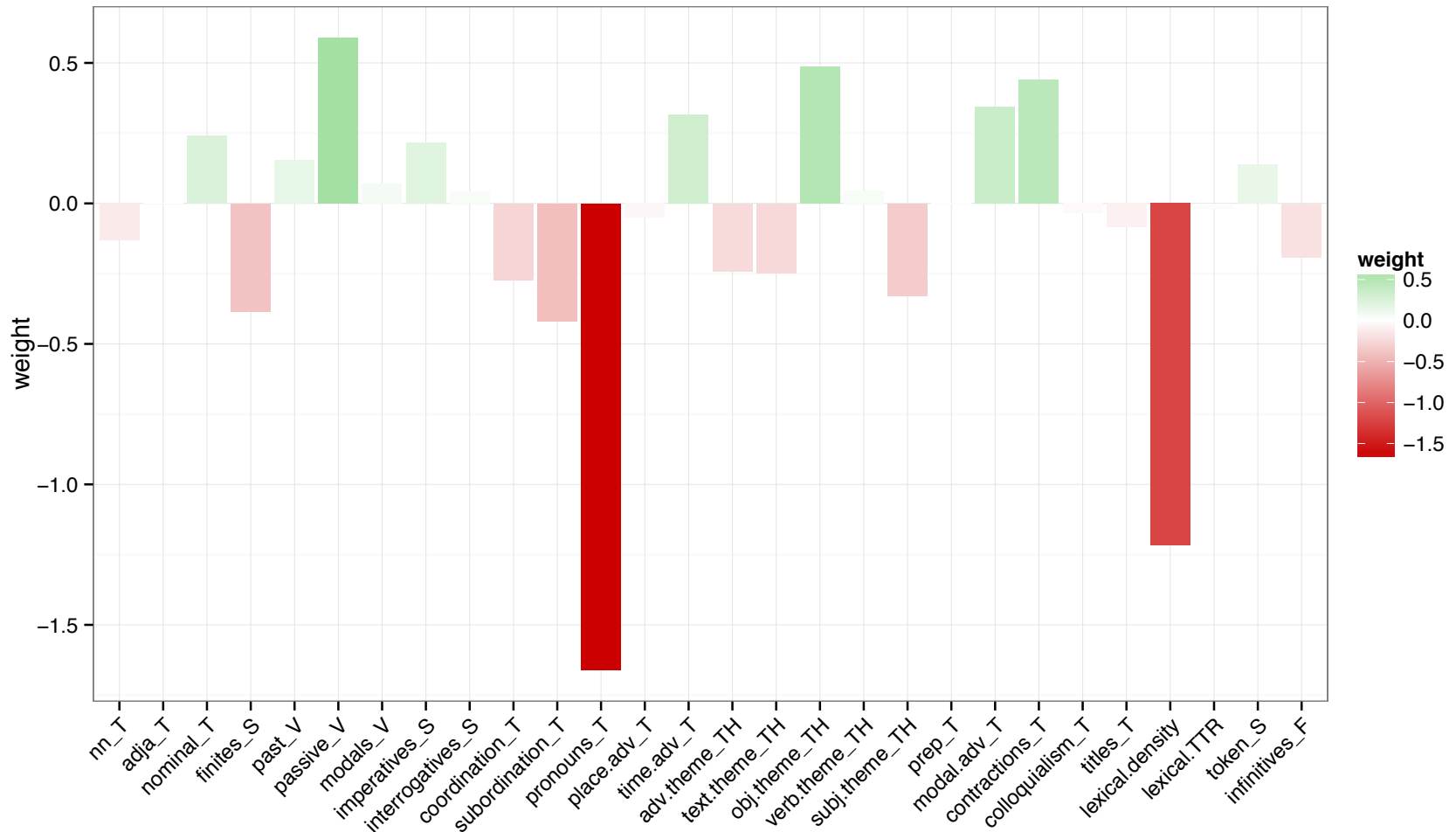
Interpreting discriminant features



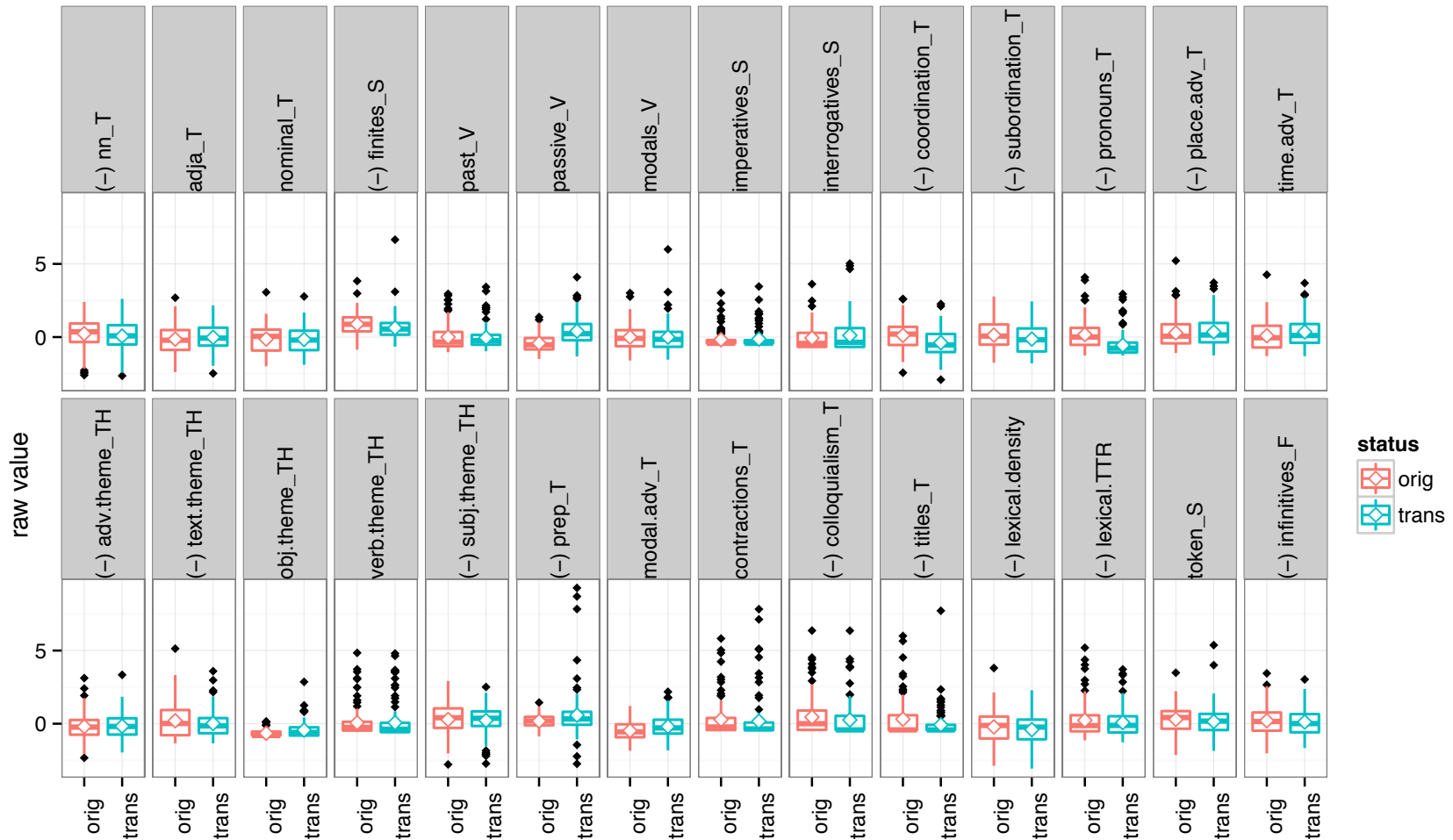
originals

translations

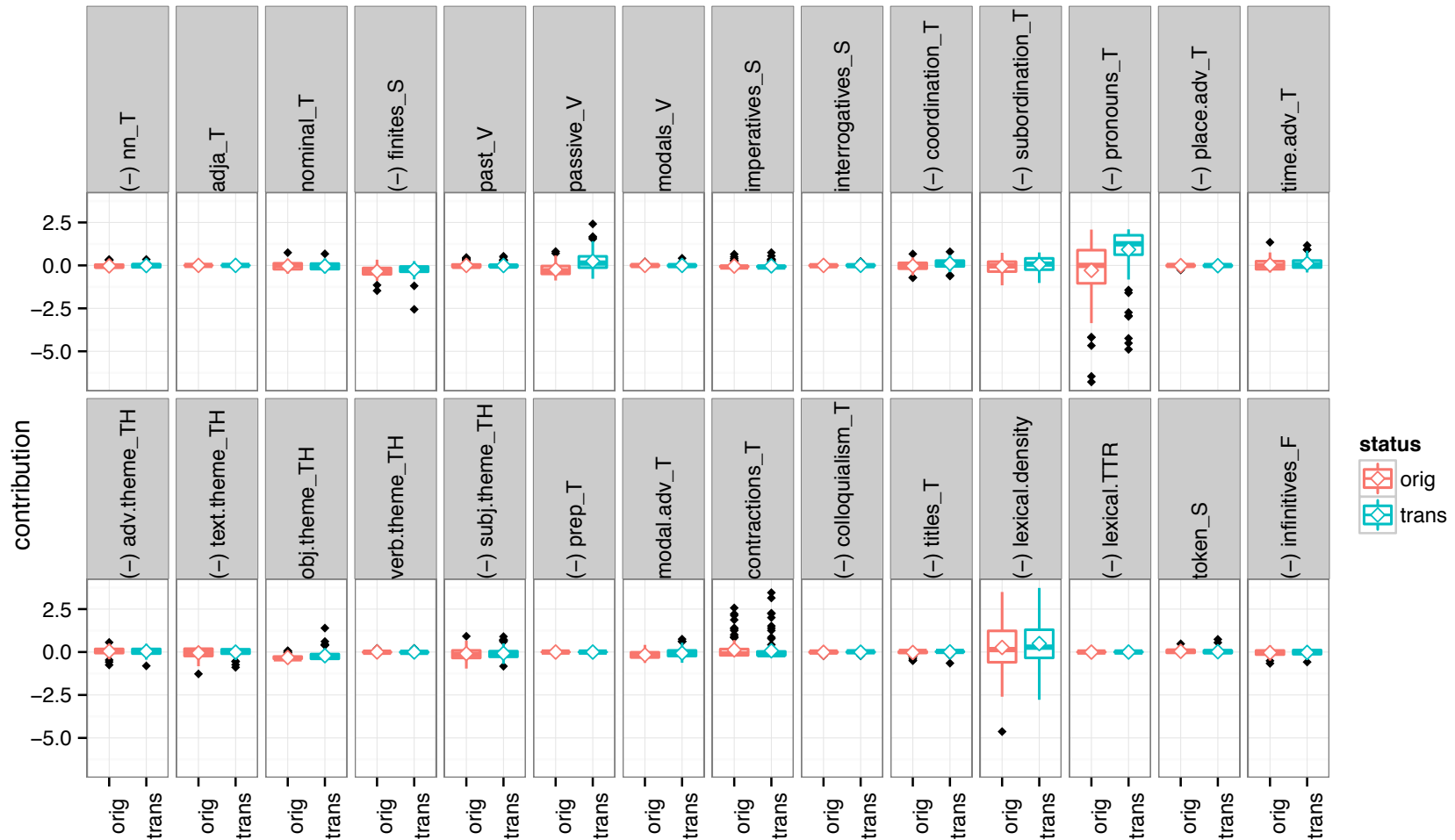
Interpreting discriminant features



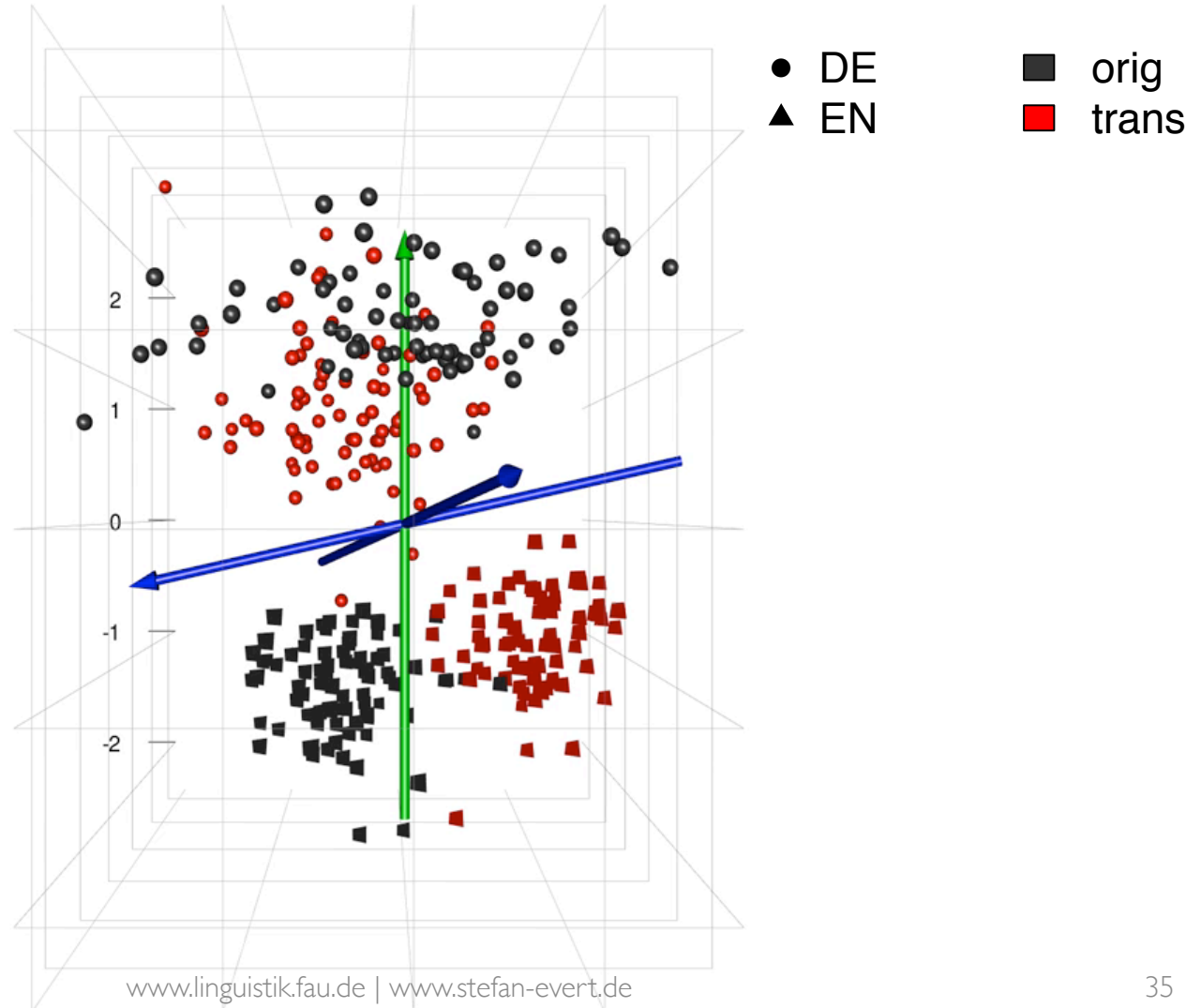
Interpreting discriminant features



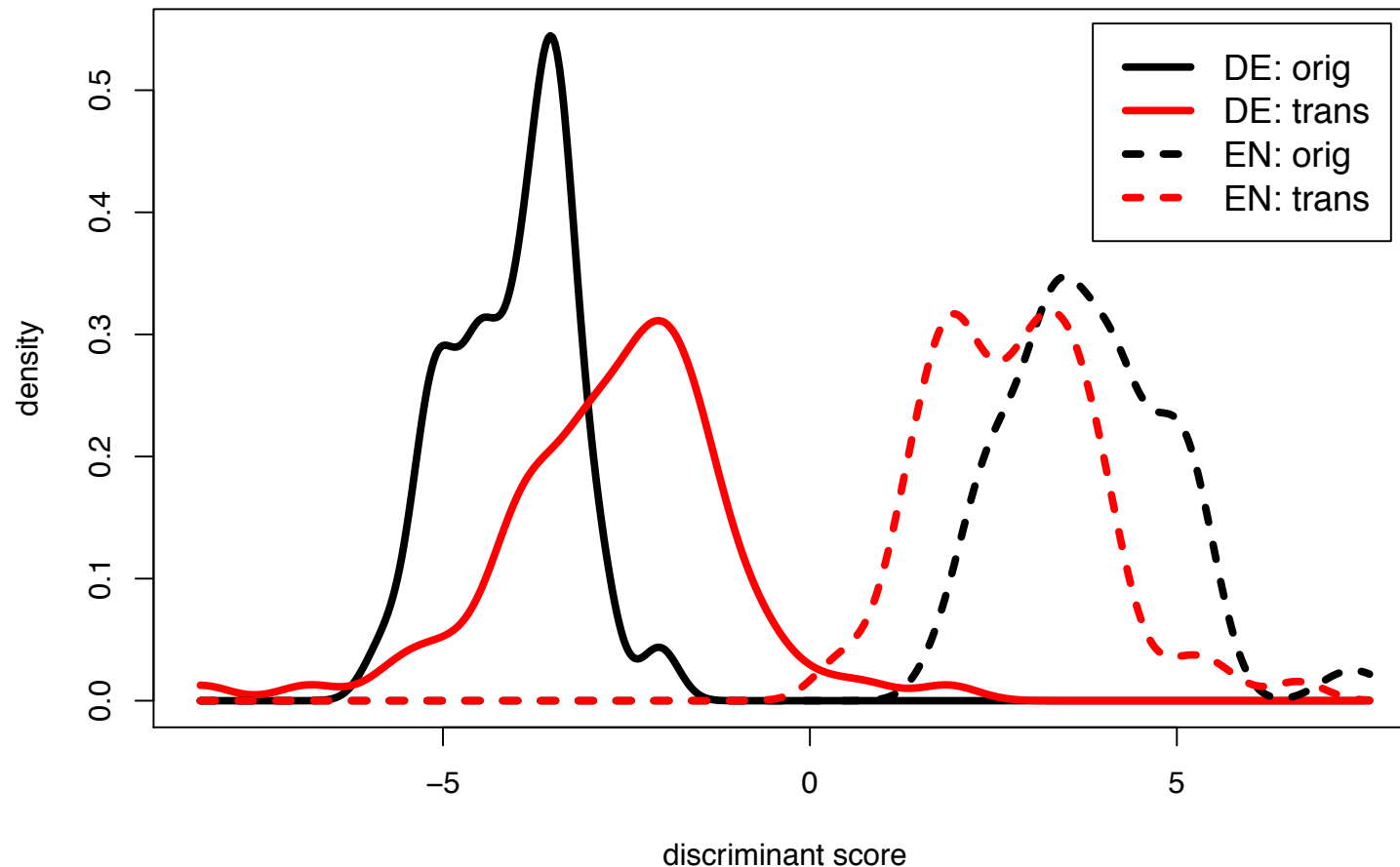
Interpreting discriminant features



Interpreting geometric configurations: German vs. English



Discriminant for DE/EN: Evidence for shining through & prestige?



Case study 2: French regional varieties

Diwersy, Evert & Neumann (2014)

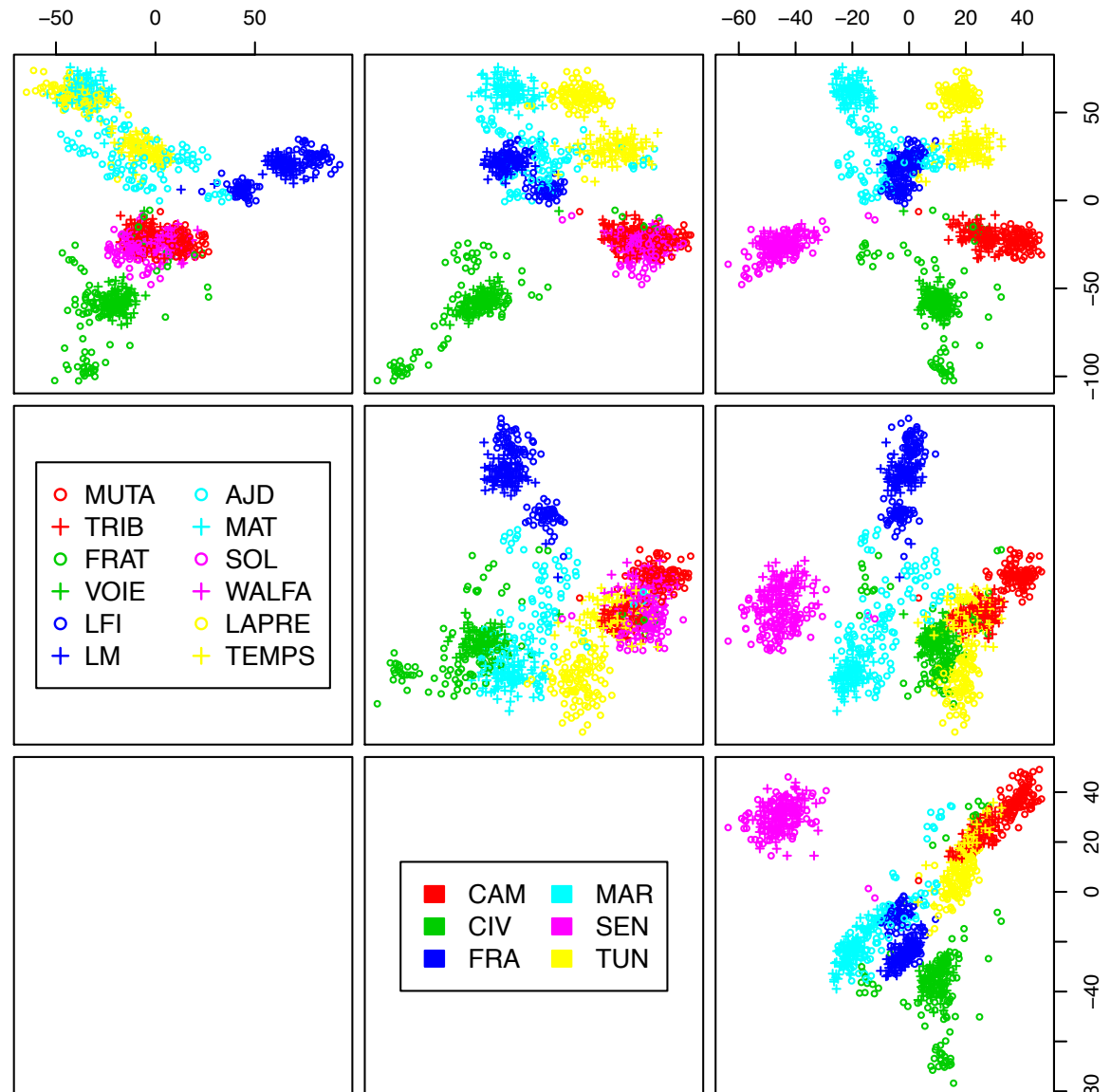
- Lexical differences in regional varieties of French
- Two nation-wide newspapers each from 6 countries
 - Cameroon, France, Ivory Coast, Morocco, Senegal, Tunisia
 - two consecutive volumes from each newspaper
 - total size approx. 14.5 million tokens
- Text samples = one week each
- Features: frequencies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

Case study 2: French regional varieties

PCA including
country-specific
words as features:
perfect separation

Design bias results
in a completely
uninteresting model

FA not applicable:
features >> texts

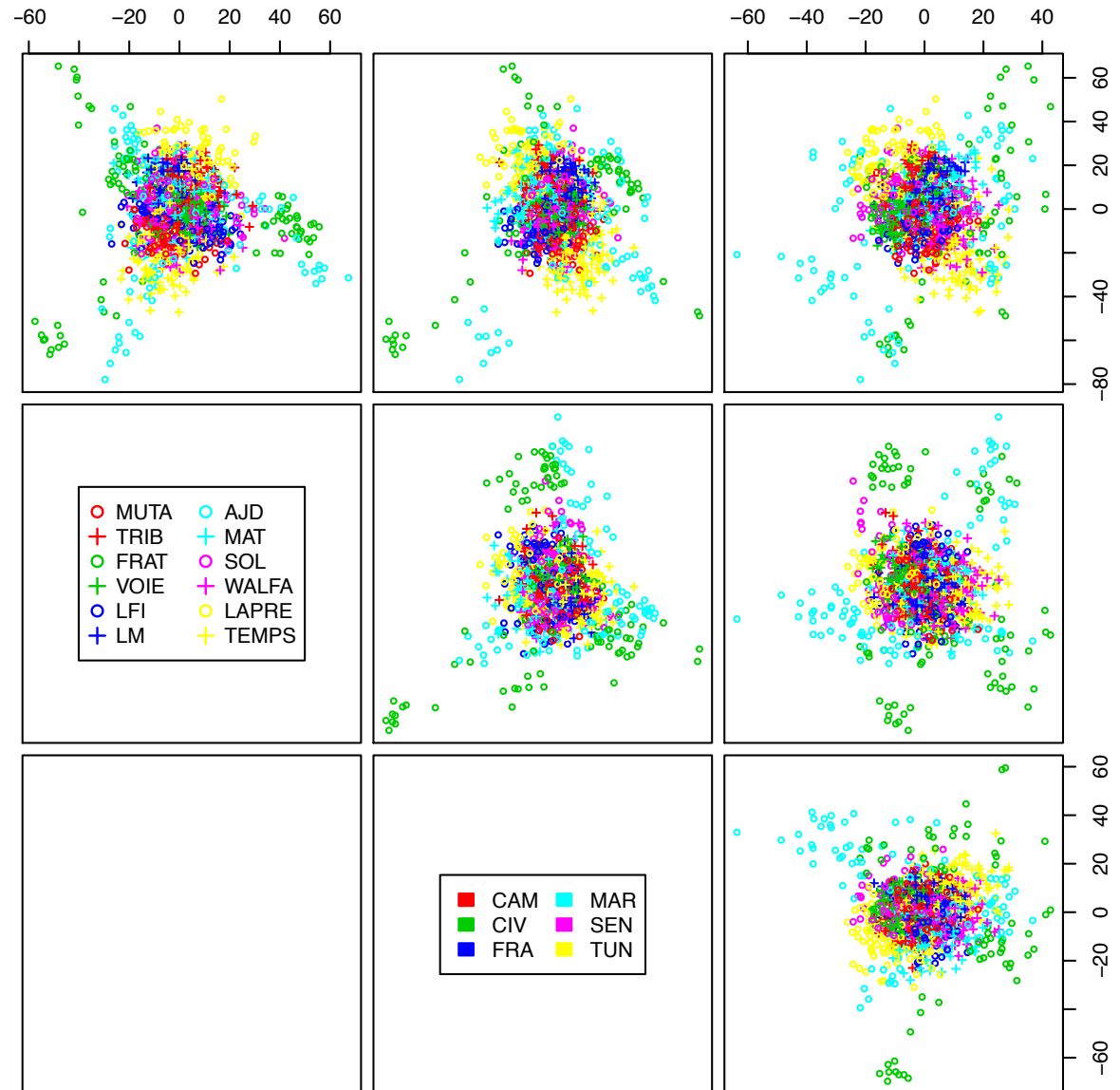


Case study 2: French regional varieties

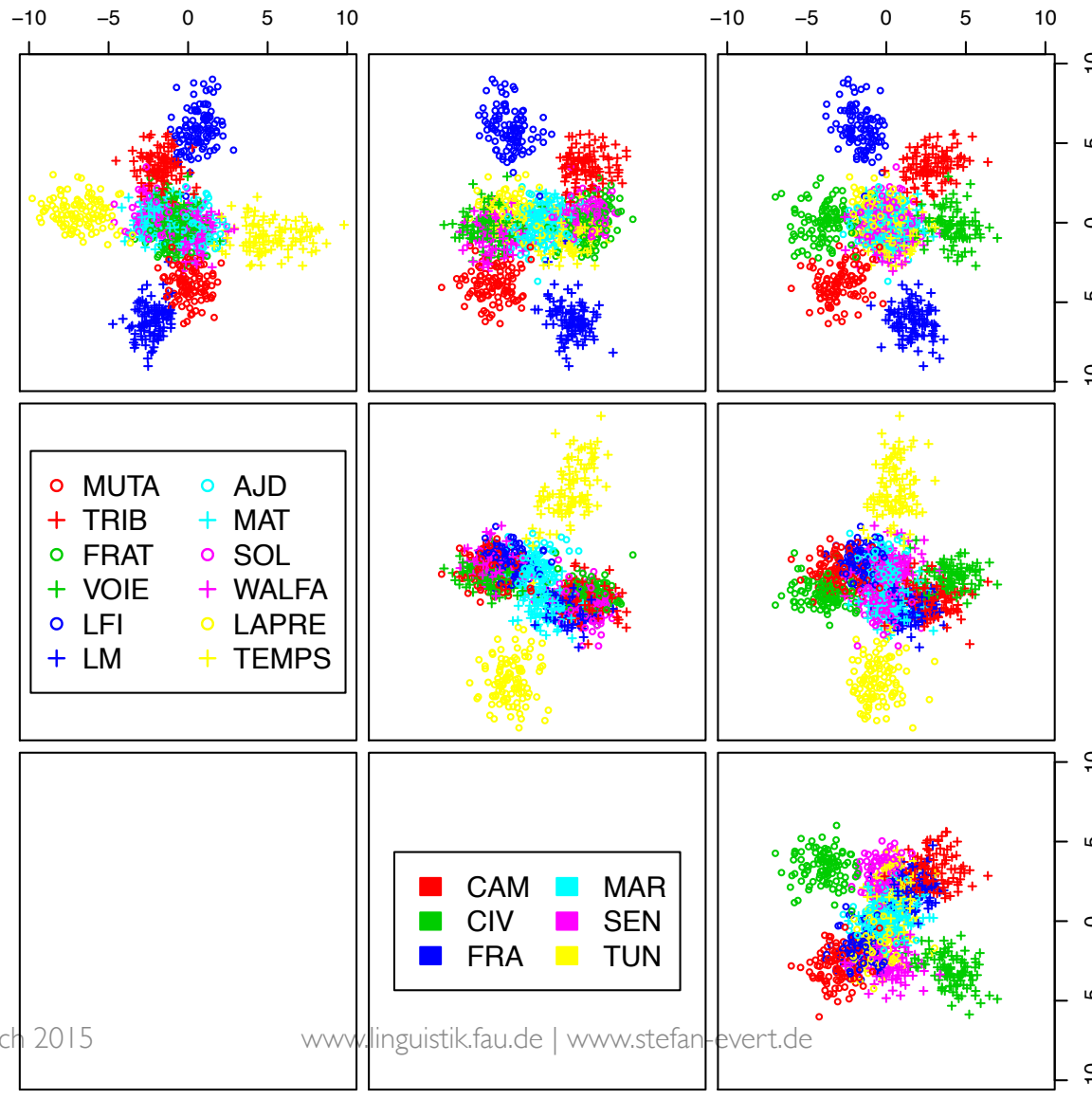
Using only shared words as features,
PCA no longer
reveals any patterns
(just a few outliers)

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

Country would presume
regional varieties exist!



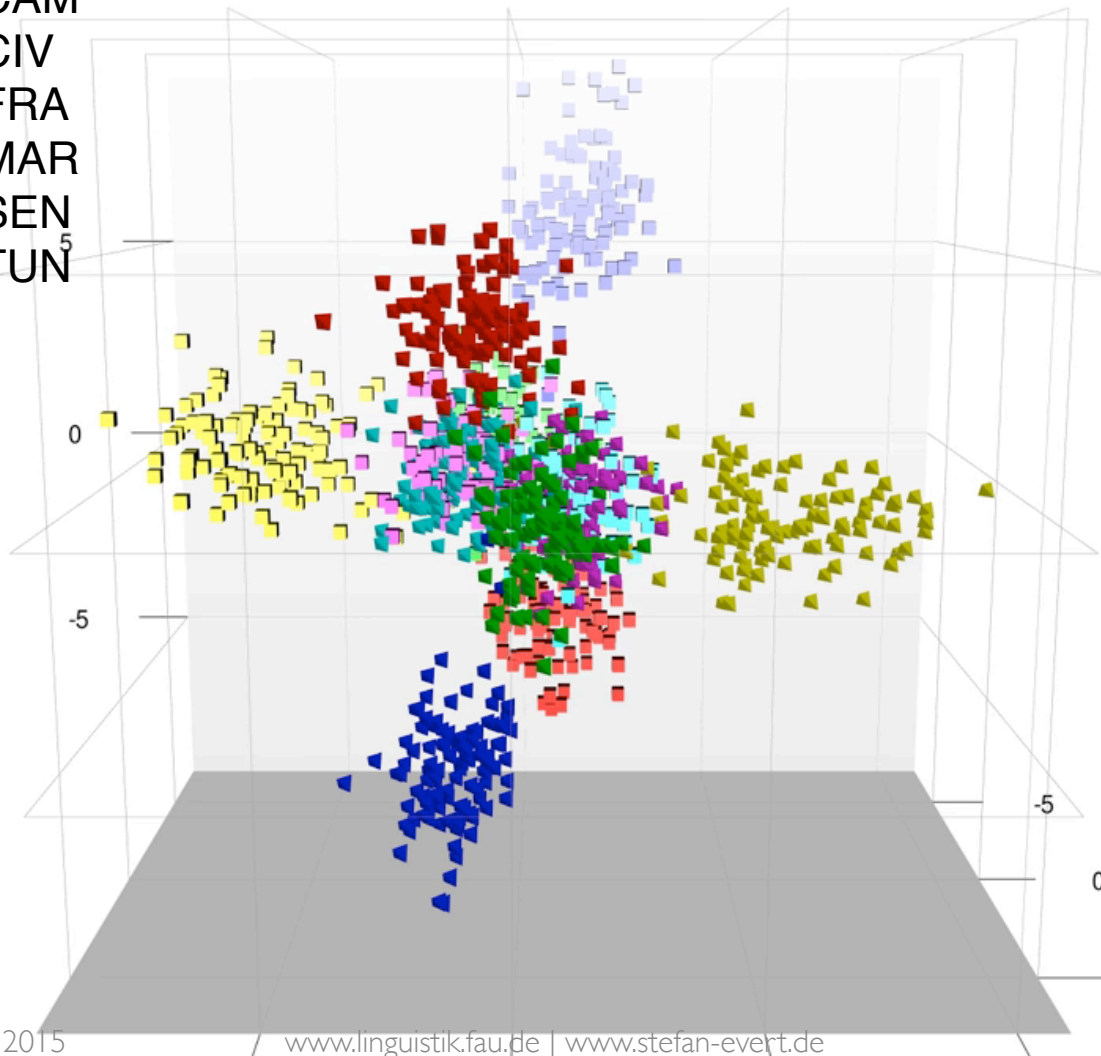
Case study 2: French regional varieties



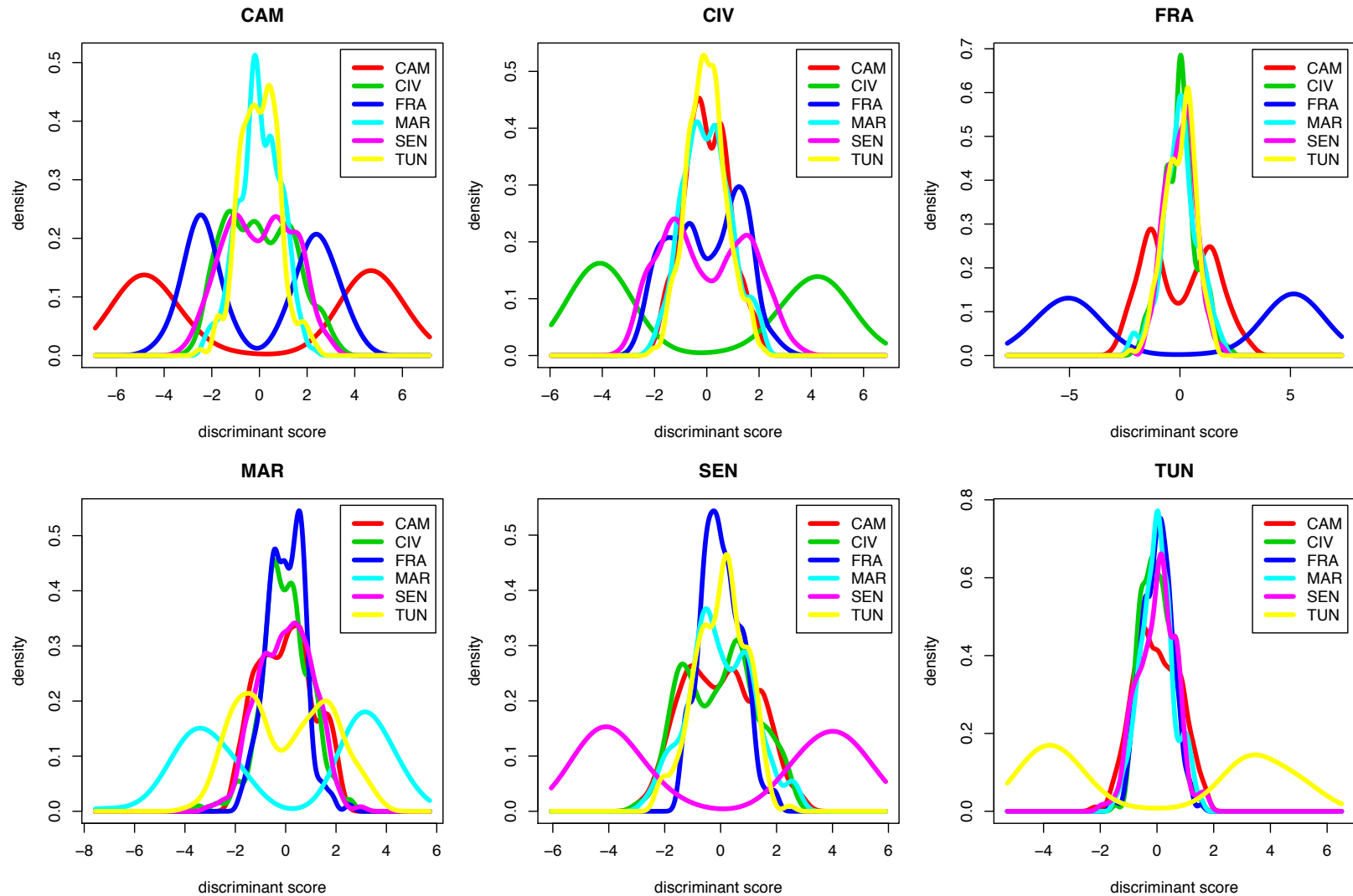
LDA dimensions (newspapers)

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

■ CAM
 ■ CIV
 ■ FRA
 ■ MAR
 ■ SEN
 ■ TUN



Discriminant axes (newspapers)





THANK YOU!

References



FRIEDRICH-ALEXANDER
UNIVERSITÄT
ERLANGEN-NÜRNBERG

PHILOSOPHISCHE FAKULTÄT
UND FACHBEREICH THEOLOGIE

- Biber, Douglas (1988). *Variation Across Speech and Writing*. Cambridge University Press, Cambridge.
- Diwersy, Sascha; Evert, Stefan; Neumann, Stella (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, Stefan & Neumann, Stella (in prep.). *The impact of translation direction on the characteristics of translated texts: a multivariate analysis for English and German*.
- Evert, Stefan; Proisl, Thomas; Schöch, Christof; Jannidis, Fotis; Pielström, Steffen; Vitt, Thorsten (2015). *Explaining Delta, or: How do distance measures for authorship attribution work? Presentation at Corpus Linguistics 2015*, Lancaster, UK.
- Gasthaus, Jan (2007). *Prototype-Based Relevance Learning for Genre Classification*. B.Sc. thesis, Universität Osnabrück, Institute of Cognitive Science.
- Koppel, Moshe; Argamon, Shlomo; Shimoni, Anat Rachel (2003). *Automatically categorizing written texts by author gender*. *Literary and Linguistic Computing*, **17**(4), 401–412.
- Neumann, Stella (2013). *Contrastive Register Variation. A Quantitative Approach to the Comparison of English and German*. de Gruyter Mouton, Berlin.