

## 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

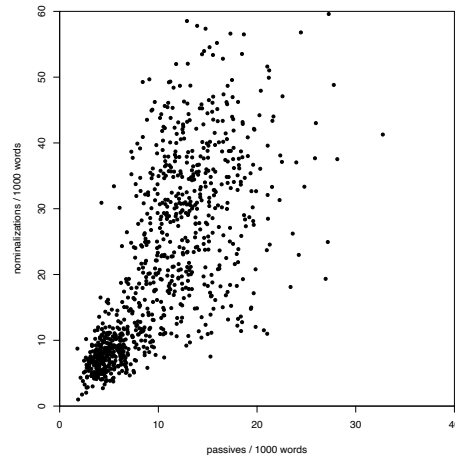
- Univariate approach
  - compare single feature across two or more conditions
  - e.g. AmE vs. BrE vs. IndE vs. ... / male vs. female / etc.
  - corpus frequency comparison
- Regression approach
  - predict single quantity from multiple explanatory factors
- Multivariate approach
  - identify common 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

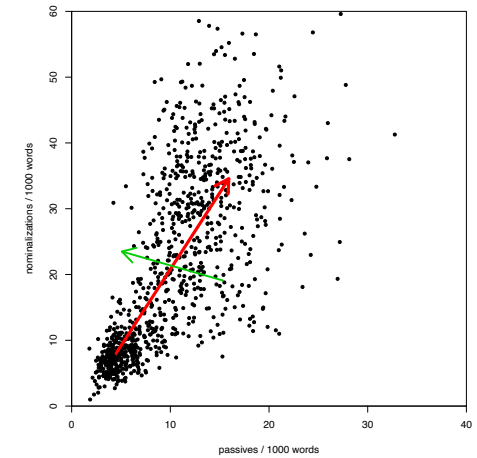
## The multivariate approach

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

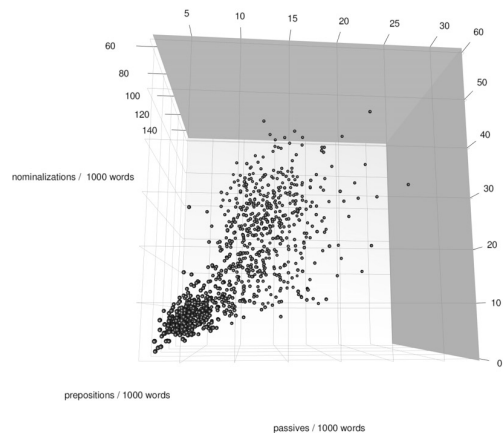


## The multivariate approach

- 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)

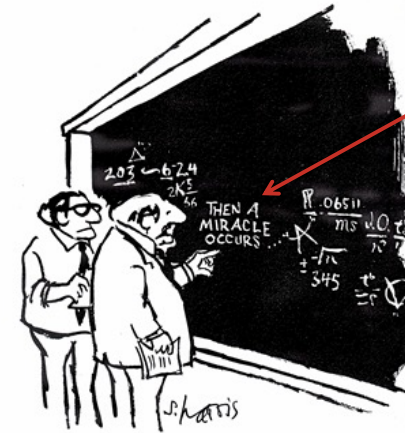
# Biber's MDA

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., above, beside, outdoors)	
5 Time adverbials (e.g., early, instantly, soon)	
C. Pronouns and pro-verbs	
6 First-person pronouns	
7 Second-person pronouns	
8 Third-person personal pronouns (excluding it)	
9 Pronoun <i>it</i>	
10 Demonstrative pronouns ( <i>that, this, these, those</i> as pronouns)	
11 Indefinite pronouns (e.g., anybody, nothing, someone)	
12 Pro-verb <i>do</i>	
D. Questions	
13 Direct WH questions	
E. Nominal forms	
14 Nominizations (ending in <i>-tion, -ment, -ness, -ity</i> )	
15 Gerunds (participle 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 <i>wh</i> -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 <i>wh</i> relatives on subject position (e.g., <i>the man who likes popcorn</i> )	
32 <i>wh</i> 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 subordinators ( <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 Down-toners (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., <i>sentence-initial 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, think, know</i> )	
57 Suggestive 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 means 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 and NOUN; ADJ; and ADJ; VERB and VERB; ADV and 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> )	



factor analysis (FA)

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

SIGIL

9

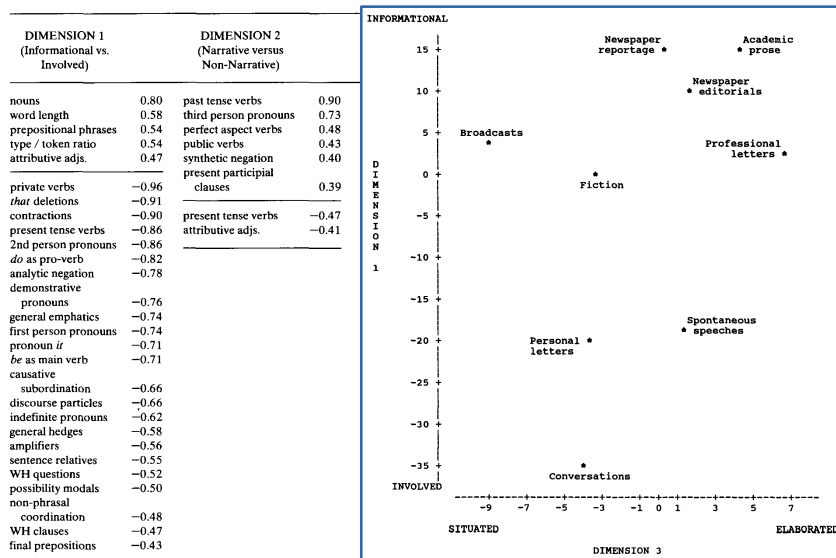
SIGIL Unit #7

www.linguistik.fau.de | www.stephanie-evert.de

10

# Biber's MDA

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



# 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?

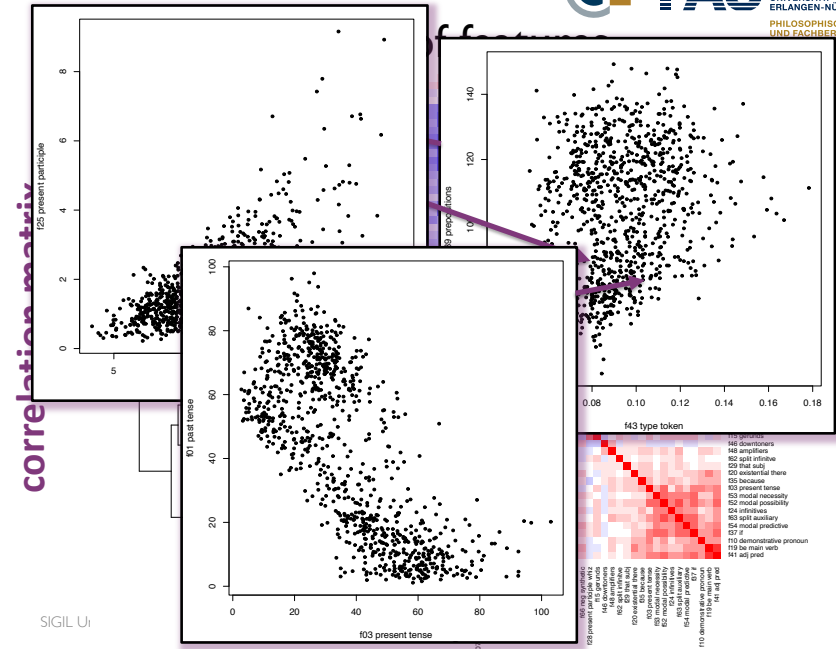
SIGIL Unit #7

www.linguistik.fau.de | www.stephanie-evert.de

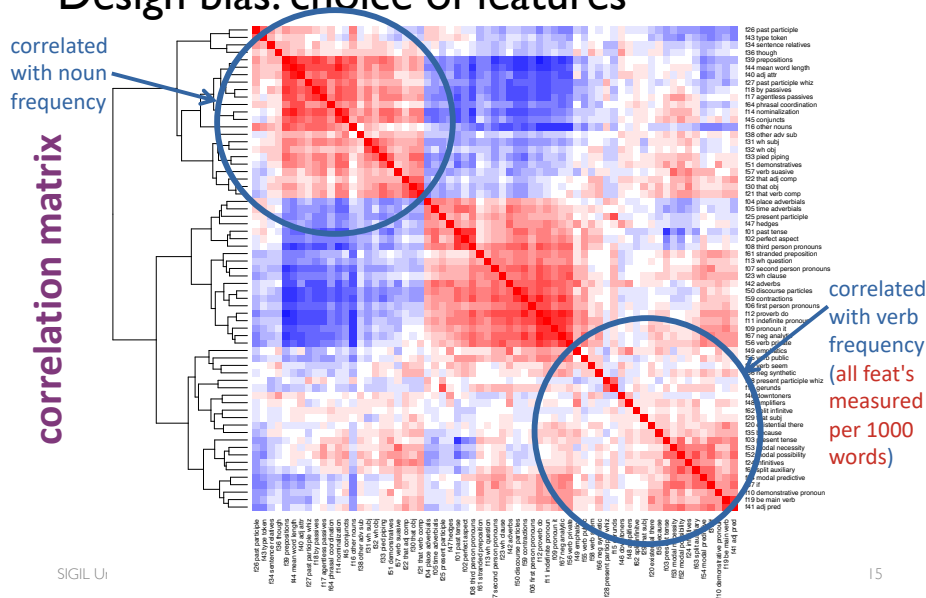
12

# Reproducing Biber's dimensions

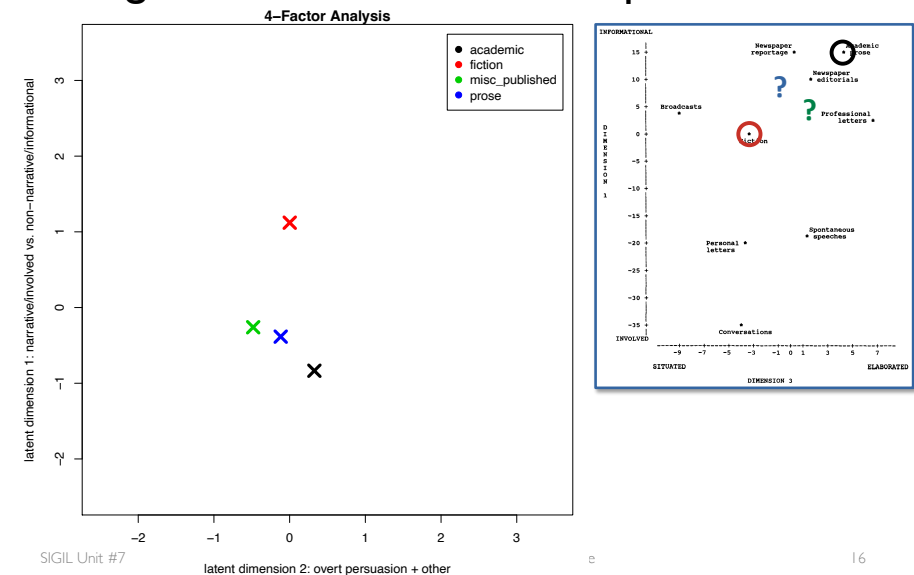
- Sample of 923 medium-length published texts from written part of British National Corpus (BNC)
- Covers 4 different text types + male/female authors
  - academic writing, non-academic prose, fiction, misc.
- Biber features extracted automatically with Python script (Gasthaus 2007)
  - all frequencies 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



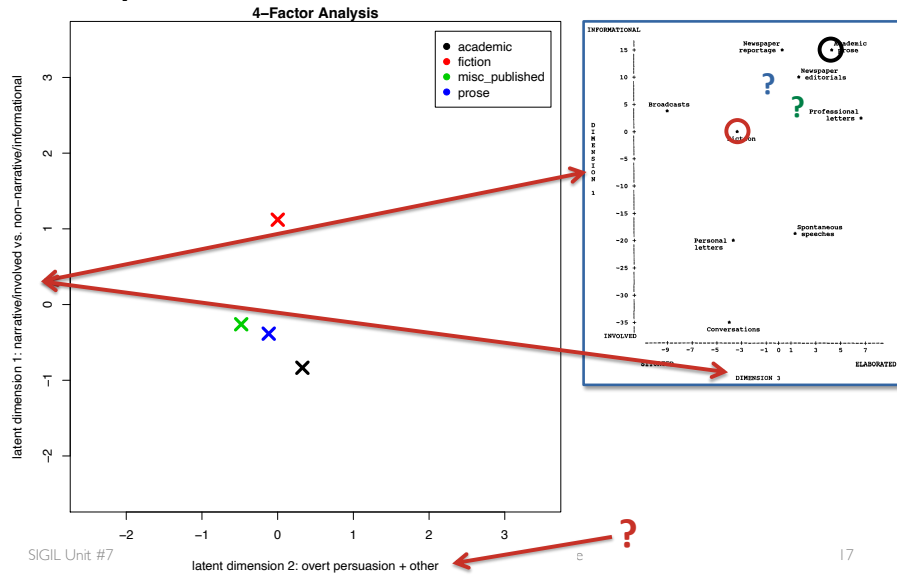
# Design bias: choice of features



# Design bias: choice of text samples



# Interpretation bias

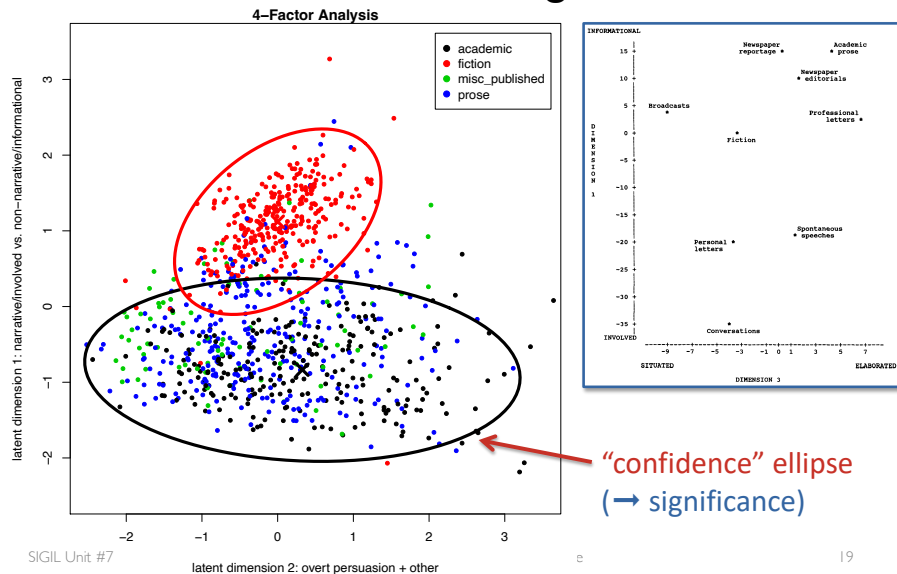


# Interpretation bias

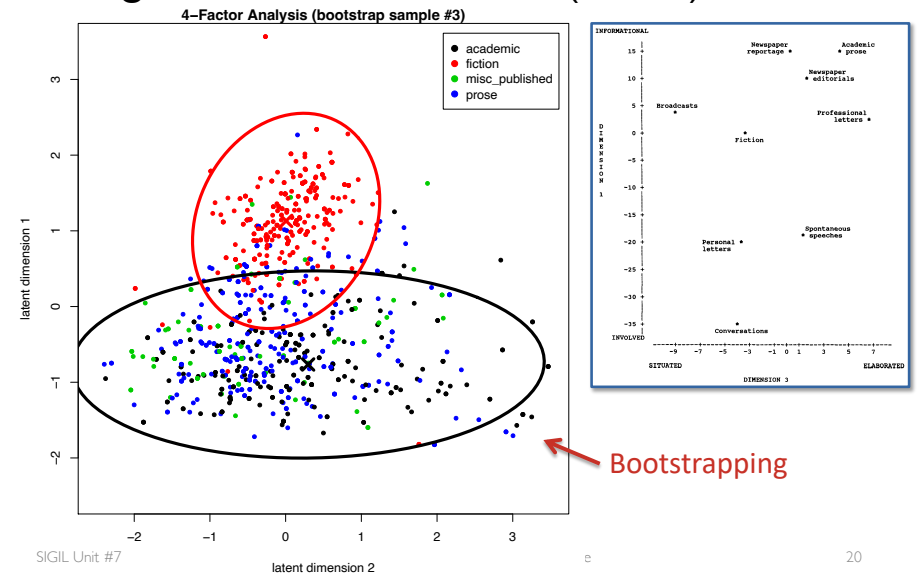
TABLE 2 Summary of the co-occurrence patterns underlying five major dimensions of English.

DIMENSION 1 (Informational vs. Involved)	DIMENSION 2 (Narrative versus Non-Narrative)	DIMENSION 3 (Elaborated vs. Situated Reference)	DIMENSION 4 (Overt Expression of Persuasion)	DIM 5 (Abs Non-A)
nouns 0.80	past tense verbs 0.90	WH relative clauses on object positions 0.63	infinitives 0.76	conjunctive
word length 0.58	third person pronouns 0.73	object positions 0.63	prediction modals 0.54	agentive
prepositional phrases 0.54	perfect aspect verbs 0.48	WH relative clauses on subject position 0.45	subordinate verbs 0.49	past participles
type / token ratio 0.54	public verbs 0.43	phrasal coordination 0.36	conditional 0.47	clausal
attributive adjs. 0.47	synthetic negation 0.40	nominalizations 0.36	subordination 0.47	BY-passive
private verbs -0.96	present participial clauses 0.39	time adverbials -0.60	necessity modals 0.46	past participles
that deletions -0.91	present tense verbs -0.47	place adverbials -0.49	split auxiliaries 0.44	WHI
contractions -0.90	attributive adjs. -0.41	other adverbs -0.46	possibility modals 0.37	other auxiliary
present tense verbs -0.86			[No complementary features]	subordinate
2nd person pronouns -0.86				[No complementary features]
do 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 -0.66				
subordination -0.66				
discourse particles -0.66				
indefinite pronouns -0.62				
copular verbs -0.59				

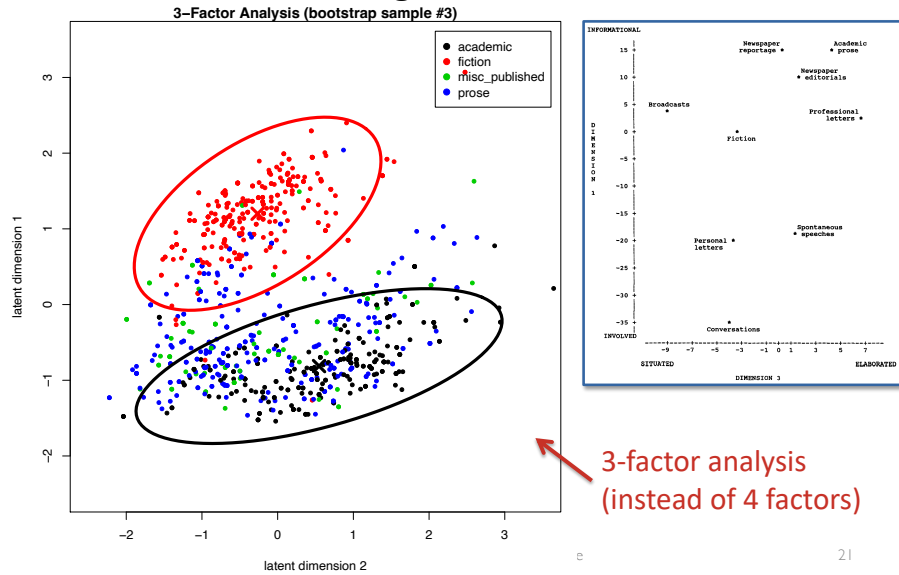
# Variation between texts is ignored



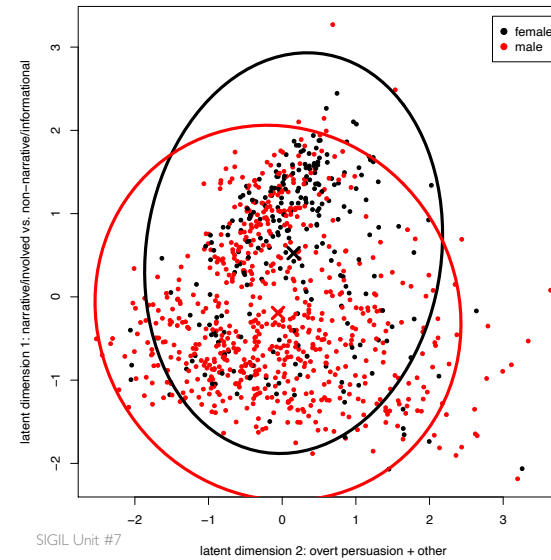
# Design bias: choice of texts (redux)



# And there's the magic number ...

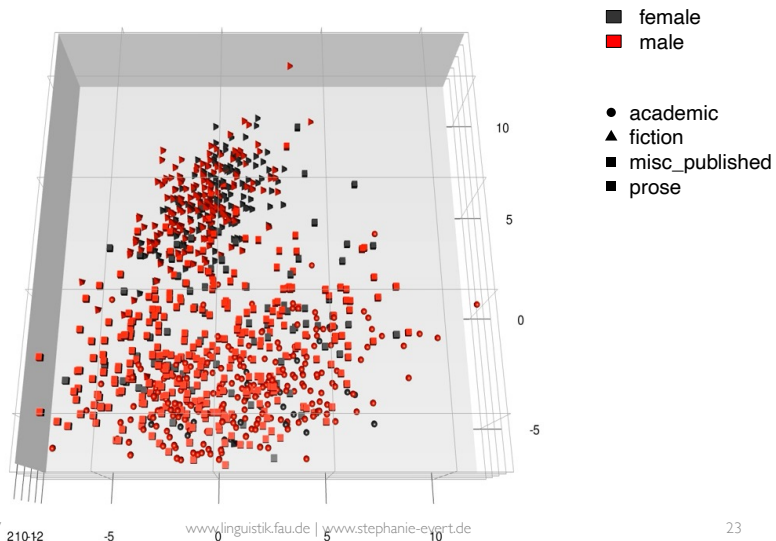


# Blindness to subtle patterns



- 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/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

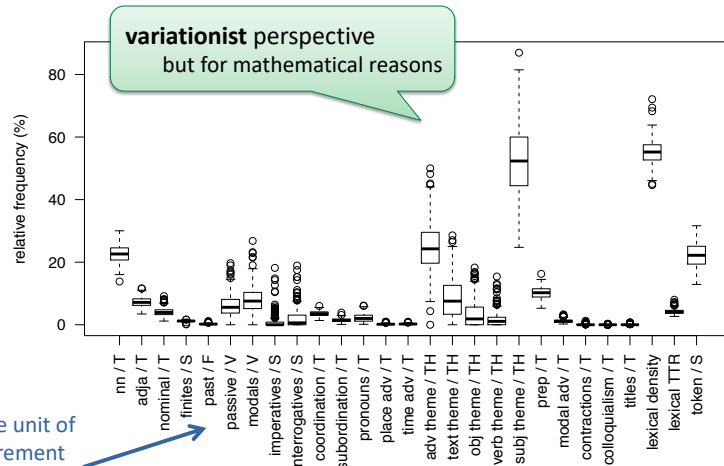
# Case study: CroCo

(Neumann 2013; Evert & Neumann 2017)

genre: language-external situation + purpose  
register: language-internal co-occurrence patterns of linguistic features

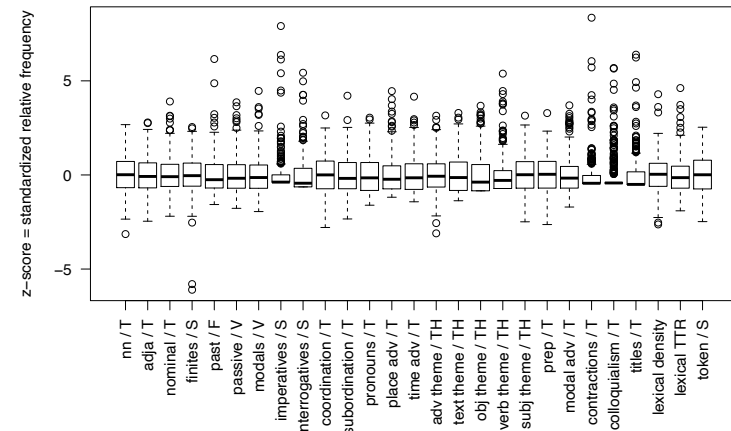
- CroCo: parallel corpus English/German
  - English-German and German-English translation pairs
  - we use 298 texts from 5 different genres (excluded: instruction manuals, tourism, fiction)
- 28 lexico-grammatical features (Neumann 2013)
  - comparable between languages
  - inspired by SFL and translation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

## Feature design: avoid “obvious” correlations

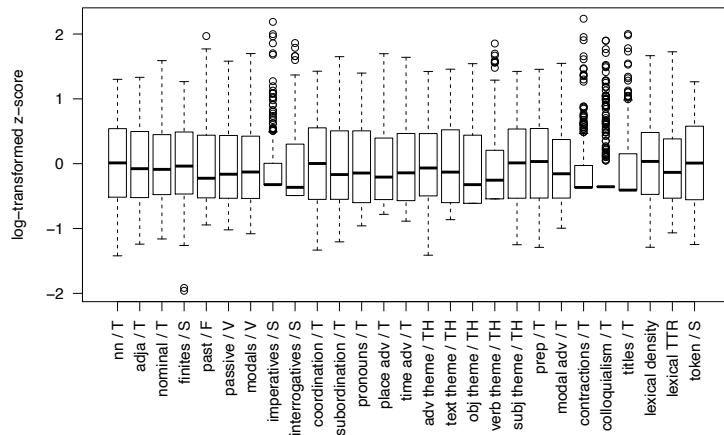


suitable unit of measurement (not always per 1000 words!)

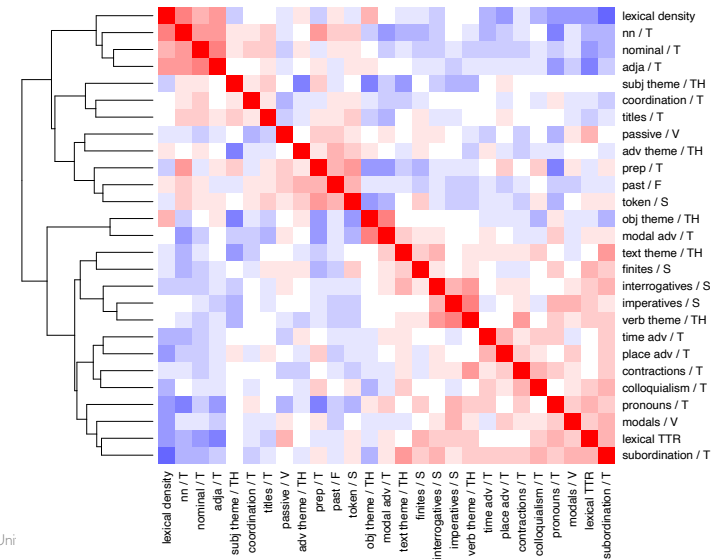
## Feature scaling: same contribution to Euclidean distances



# Feature scaling: optional signed log transformation



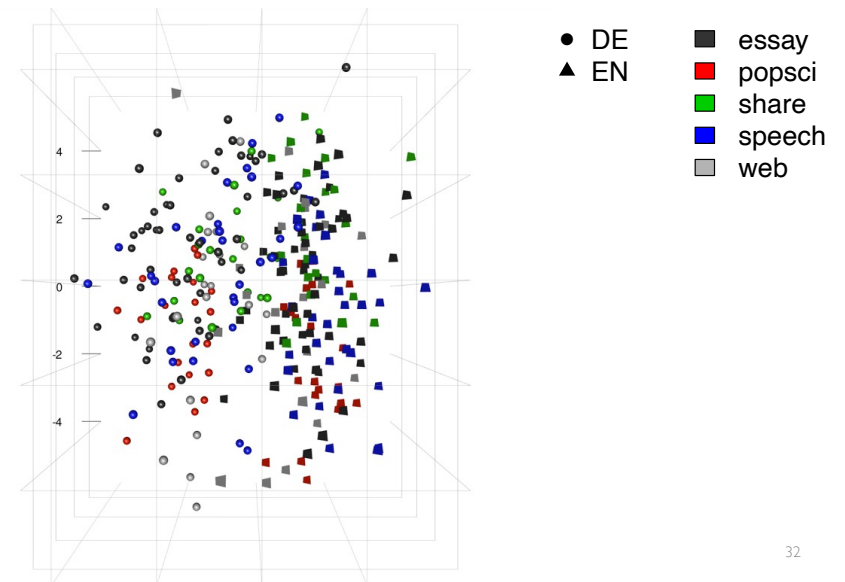
# CroCo: correlation matrix



# Latent dimensions as perspective on data configuration

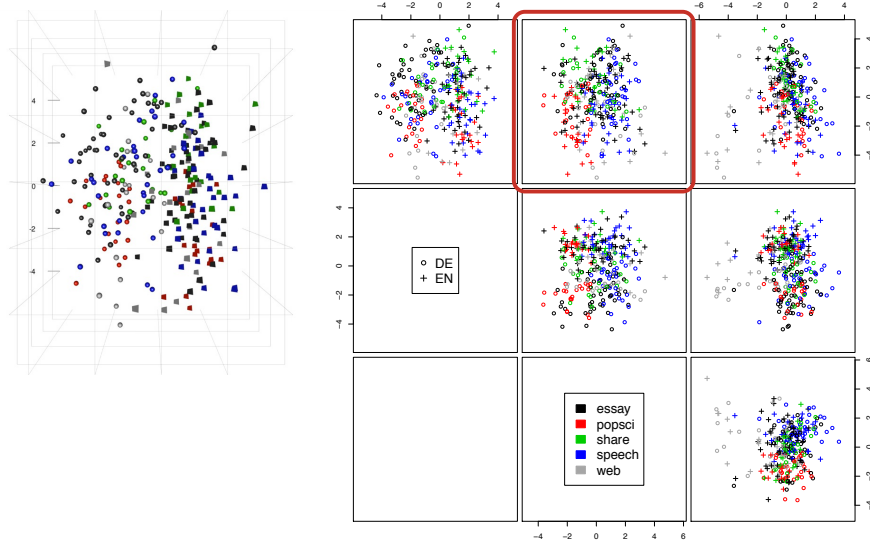
- Instead of “magical” latent dimensions we focus on **orthogonal projections** as perspectives on the data
  - cf. photograph as 2D perspective on 3D object
- Different perspectives highlight different aspects
- Multivariate analysis → choice of perspective
  - **principal component analysis (PCA)** = 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



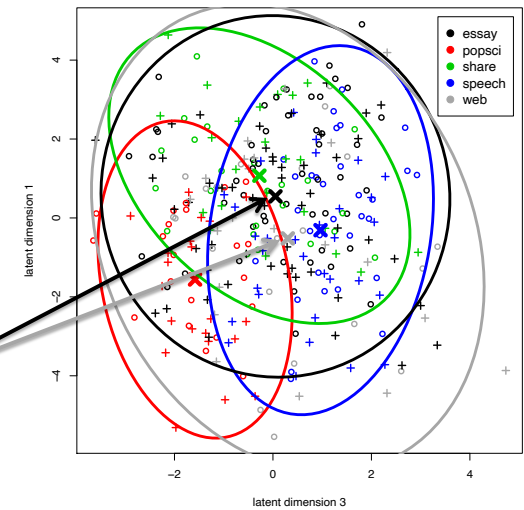


## CroCo: 4-dimensional projection



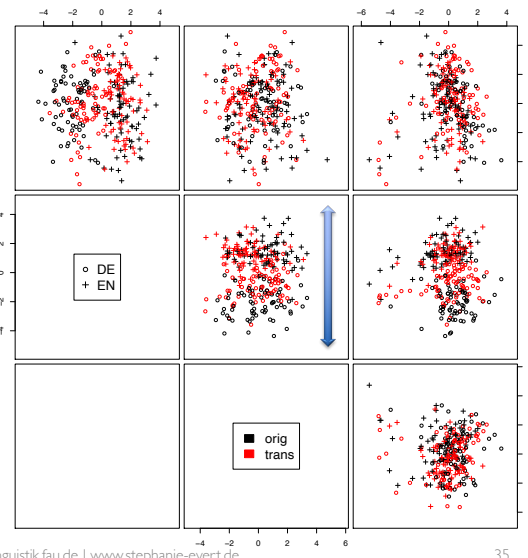
## CroCo: genre distribution

- 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?

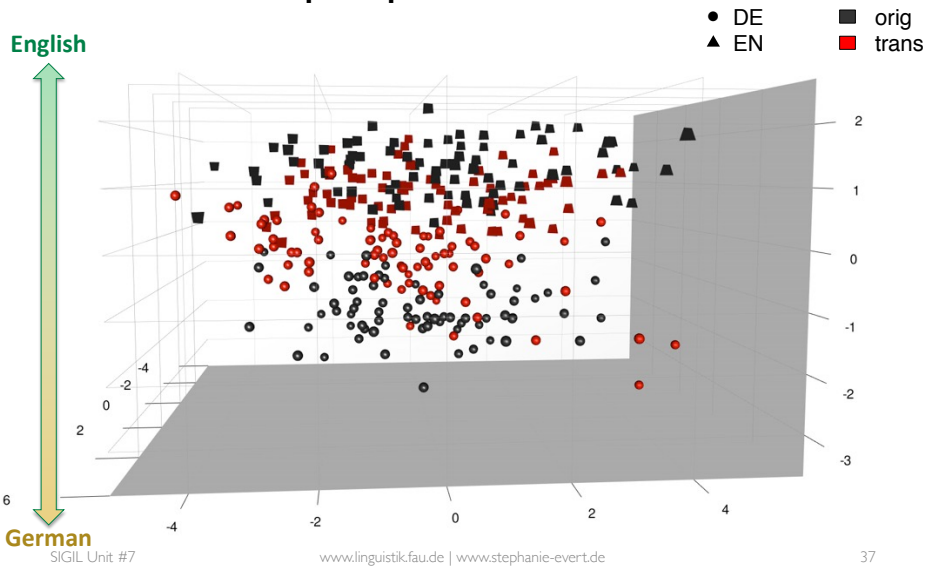
- PCA dim's can't separate translations from original texts
  - 62.1% accuracy on first 3 PCA dim's
- But SVM machine learner can do this with >80% accuracy
  - RBF kernel
  - 10-fold c.v.
- Hints at **shining through**, but no clear-cut evidence



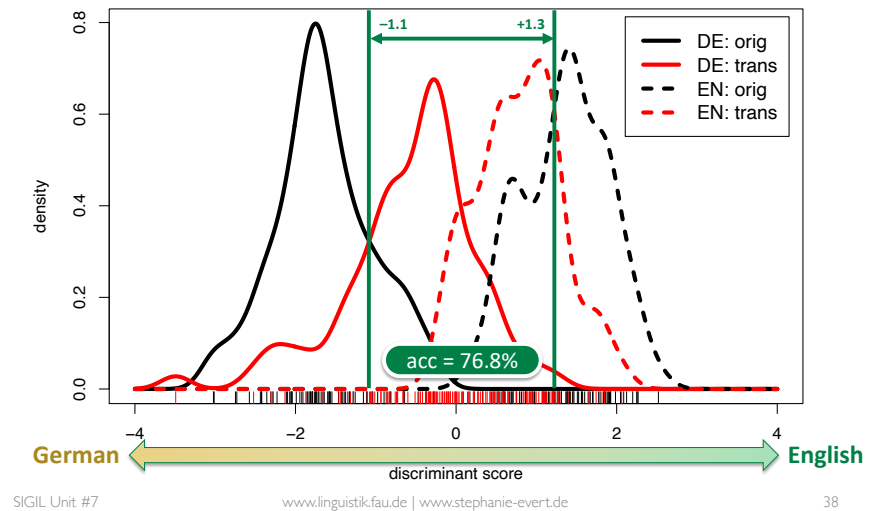
## Minimally supervised LDA

- 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

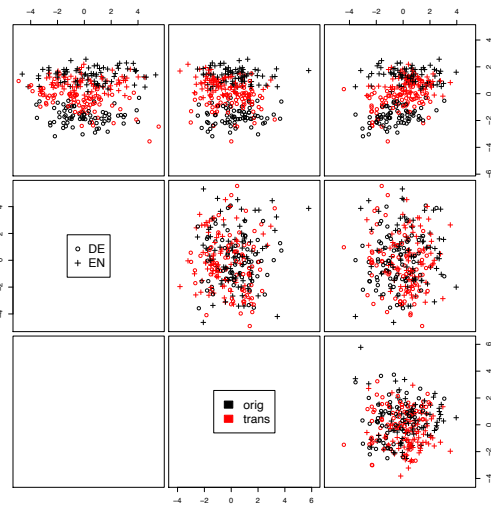


# Discriminant for DE vs. EN confirms shining through & prestige effect

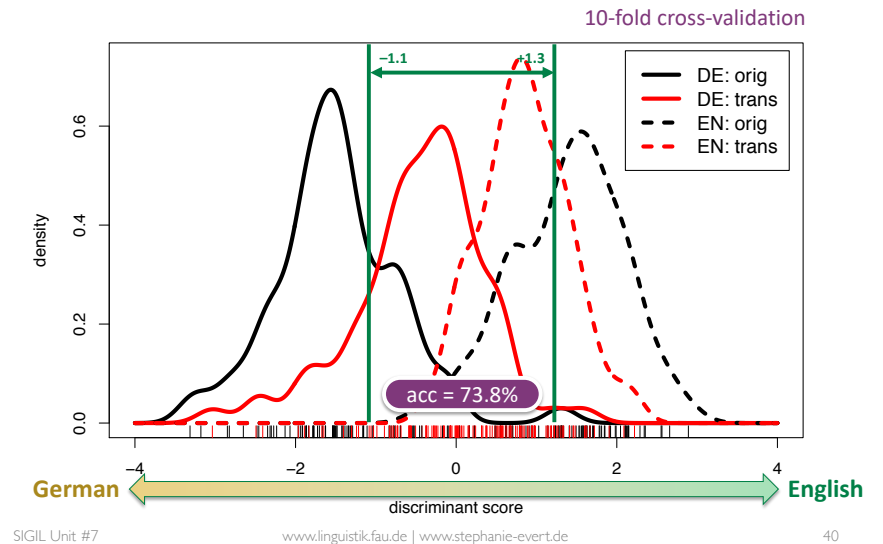


# LDA significance: bootstrapping / cross-validation

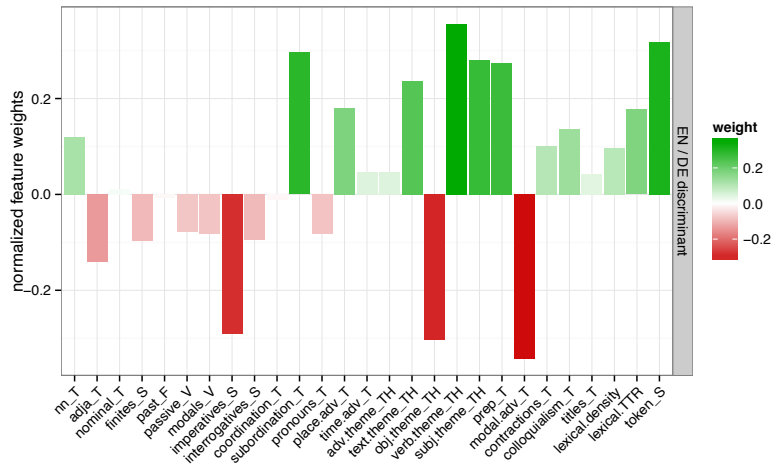
- LDA is a supervised ML technique → overtrained?
  - Would we find the same discriminant if we trained on a different set of texts?
- Verification with **bootstrap resampling** or **10-fold cross-validation**
  - LDA trained on 90% of data
  - discriminant axis shows “wobble” of approx. 10°
- Discriminant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



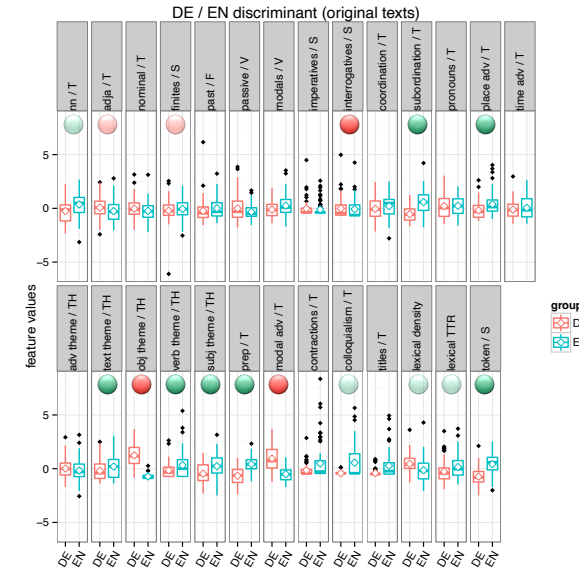
# Cross-validated discriminant



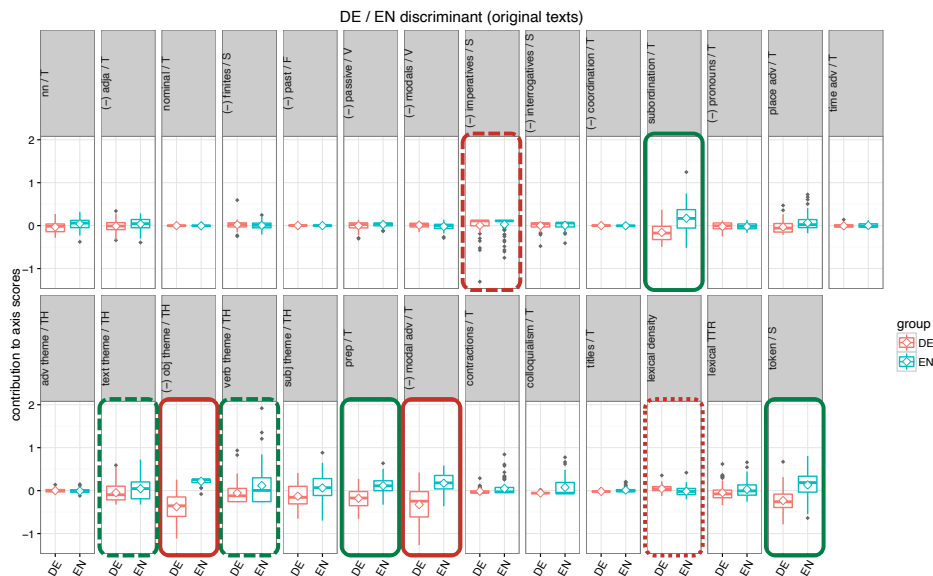
# Interpreting discriminant features



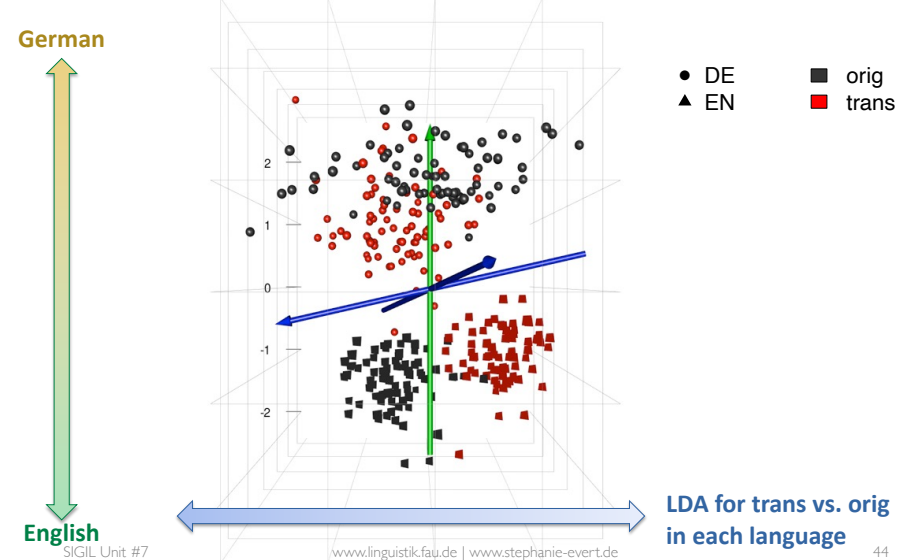
# Interpreting discriminant features



# Interpreting discriminant features



# Unravelling translationese



## Case study 2: French regional varieties

(Diwersy, Evert & Neumann 2014)

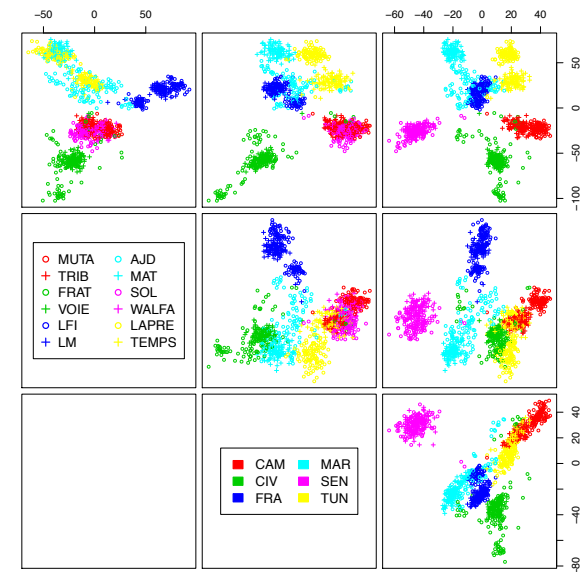
- Lexical differences in regional varieties of French
- Two nation-wide newspapers each from 6 countries
  - Cameroon, France, Ivory Coast, Morocco, Senegal, Tunisia
  - two consecutive volumes from each newspaper
  - total size approx. 14.5 million tokens
- Text samples = one week each
- Features: frequencies 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  $\gg$  texts

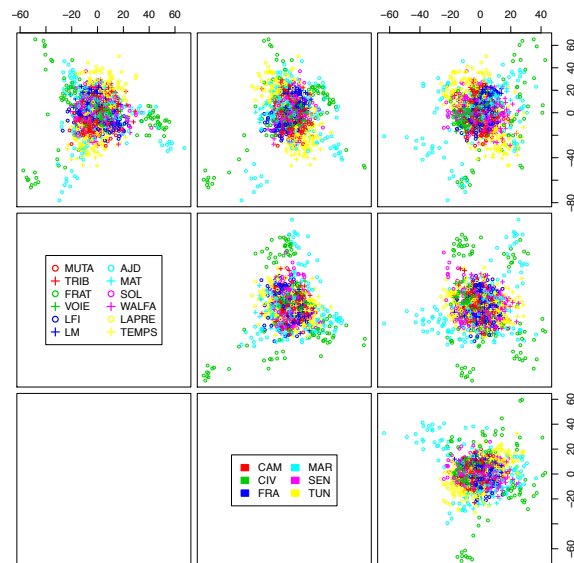


## FRV: PCA dimensions

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

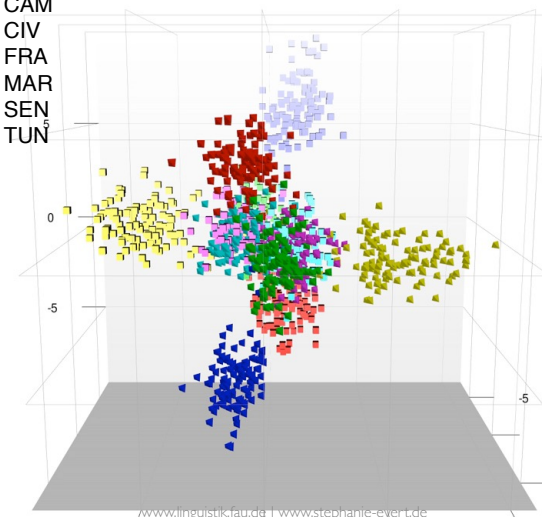
Use LDA to find a meaningful perspective, based on newspaper source

Country would presume regional varieties exist!

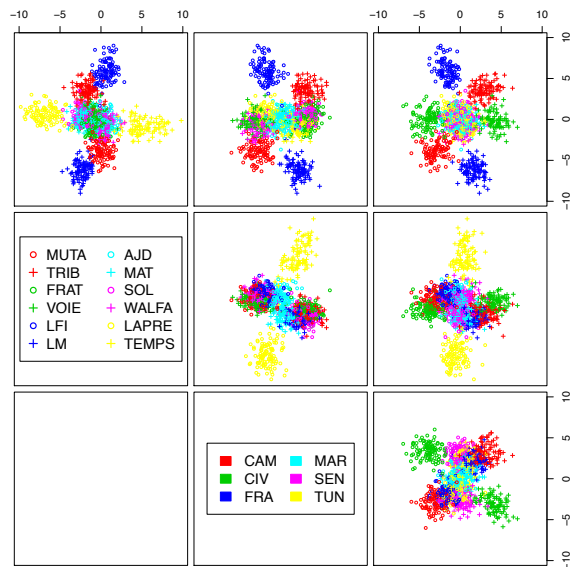


## FRV: LDA dimensions (newspapers)

- MUTA (red square)
- TRIB (red triangle)
- FRAT (green square)
- VOIE (green triangle)
- LFI (blue square)
- LM (blue triangle)
- AJD (cyan square)
- MAT (cyan triangle)
- SOL (magenta square)
- WALFA (magenta triangle)
- LAPRE (yellow square)
- TEMPS (yellow triangle)
- CAM (red square)
- CIV (green square)
- FRA (blue square)
- MAR (cyan square)
- SEN (magenta square)
- TUN (yellow square)



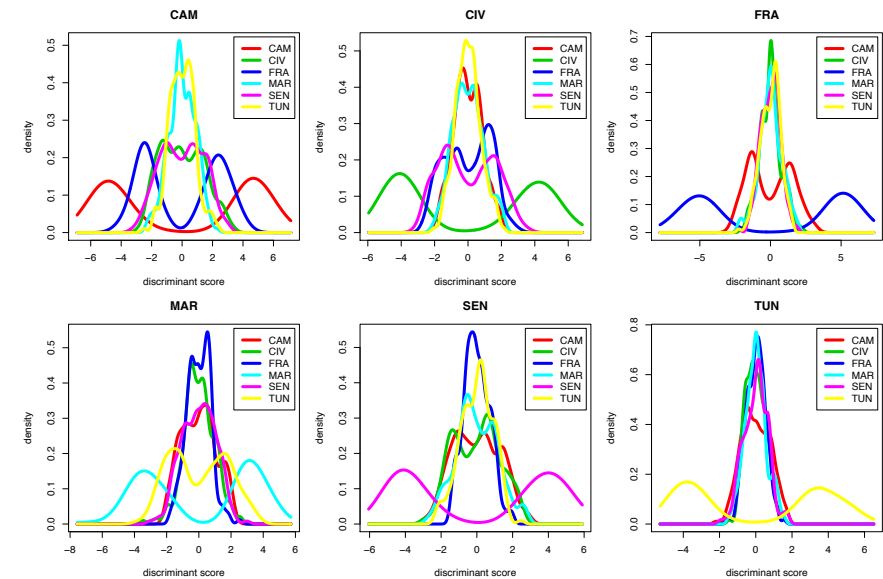
## FRV: LDA dimensions (newspapers)



SIGIL Unit #7

49

## 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.; Shimoni, 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.

SIGIL Unit #7

www.linguistik.fau.de | www.stephanie-evert.de

51