

Morpho-Syntactic Properties of Kiezdeutsch: Lexicalised POS Regression Analyses

[Poster 118](#)

[Diego Frassinelli](#) (University of Konstanz);
Gabriella Lapesa, Reem Alatrash, Dominik Schlechtweg,
Sabine Schulte im Walde (University of Stuttgart)

Our Work

We propose a quantitative **empirical framework** for larger-scale bottom-up analyses of the **lexical** and **morpho-syntactic** properties of language varieties

Case study: KiezDeutsch vs. Standard German

Results:

- Large scale analysis of linguistic properties
- Consistent with the predictions drawn from the theoretical literature
- Successful method to identify the most salient properties that distinguish one language variety from another

Kiezdeutsch: an independent variety of German

Spoken by teenagers from multi-ethnic urban neighborhoods in casual conversations among peers → **interesting morpho-syntactic properties:**

1. Bare nouns lacking determiners:
 - Hast du **ein** Problem? (Do you have problem?)
2. Bare nouns lacking prepositions:
 - Ich gehe **ins** Kino. (I go cinema)
3. Lack of copula verbs:
 - Er **ist** aus Kreuzberg. (He from Kreuzberg.)
4. Verb first declaratives:
 - **Ich** Wollte **ich** keine Hektik machen da drinne. (I didn't want to make any hectic in there.)
5. Subject-Verb-Object word order in sentences beginning with an adverb:
 - Jetzt **bin** ich **bin** 18. (Now I am 18.)

Materials

KiDKo

The KiezDeutsch Korpus*

Spoken

Casual everyday conversations in Berlin

Teenagers (14 - 17 yo), self-recording

Daily life topics

63K sentences, 359K tokens

POS tagged: Stuttgart-Tübingen tagset

Very short sentences: 8.8 tokens/sentence

* *Multi-ethnic subcorpus*

GRAIN

The German RAdio INterviews corpus

Spoken

Interviews on the public radio

Professionals (adults), recorded

Social and political topics

14K sentences, 221K tokens

POS tagged: Stuttgart-Tübingen tagset

Longer sentences: 26.7 tokens/sentence

Method: Z-score Analysis

- Logistic Regression models

corpus type_{KidKo/GRAIN} ~ POS_{0/1}

- One POS at the time: NOUN_{0/1}, VERB_{0/1}, ...
- Analysis of the z-scores:
 - **Sign** for the direction of the effect (KidKo positive, GRAIN negative)
 - **Value** for the level of uncertainty of the prediction (bigger is more predictive)
 - **Significance** test (alpha corrected p-value < 0.001)

Study 1a: Unigram POS Analysis

10 logistic regression models, one for each coarse-grain POS

Results for KiezDeutsch:

1. Significantly less determiners
 - Hast du ~~ein~~ Problem? (Do you have problem?)
2. Significantly less adpositions (prepositions)
 - Ich gehe ~~ins~~ Kino. (I go cinema)
3. More pronouns and verbs than nouns
 - Short sentences and a lot of direct speech

GRAIN		KidKo	
POS	z-sc.	POS	z-sc.
DET	-54.27	PRT	59.32
NOUN	-44.10	PRON	37.50
ADP	-38.45	ADV	22.10
CONJ	-18.17	VERB	21.92
ADJ	-11.67	NUM	7.30

Study 1b: Lexicalized Unigram POS Analysis

Nouns

GRAIN	trans.	z-score	KidKo	trans.	z-score
Menschen	<i>humans</i>	-8.23	Alter	<i>age</i>	11.19
Frage	<i>question</i>	-6.23	Schule	<i>school</i>	8.11
Thema	<i>topic</i>	-6.20	Euro	<i>euro</i>	7.85
Land	<i>country</i>	-6.02	Stunden	<i>hours</i>	7.33
Herr	<i>Mr.</i>	-5.99	Spiel	<i>game</i>	7.23
Prozent	<i>percent</i>	-5.65	Hause	<i>home</i>	7.17
Europa	<i>Europe</i>	-5.31	Ahnung	<i>idea</i>	6.88
Jahren	<i>years</i>	-4.70	Spaß	<i>fun</i>	6.72
Gesellschaft	<i>society</i>	-4.29	Minuten	<i>minutes</i>	6.34
Ende	<i>end</i>	-4.07	Mal	<i>times</i>	6.20

1. Slang forms (Alter)
2. Daily topics (school, fun)

Verbs

GRAIN	trans.	z-score	KidKo	trans.	z-score
habe	<i>have</i>	20.13	werden	<i>will be</i>	-16.81
war	<i>was</i>	11.72	haben	<i>have</i>	-13.68
weiß	<i>know</i>	10.84	wird	<i>will</i>	-12.61
gesehen	<i>seen</i>	6.58	sind	<i>are</i>	-12.35
warte	<i>wait</i>	6.27	müssen	<i>must</i>	-12.27
mache	<i>make</i>	6.18	gibt	<i>give</i>	-11.34
gesagt	<i>said</i>	6.18	können	<i>can</i>	-9.30
bin	<i>am</i>	6.09	wollen	<i>want</i>	-8.69
mach	<i>make</i>	6.05	brauchen	<i>need</i>	-7.08
gemacht	<i>made</i>	5.90	sagen	<i>say</i>	-7.03

1. Needing, having, obligations
2. Simpler structures

Study 2: Trigram POS Analysis

1,245 logistic regression models, one for each POS triplet

GRAIN		KidKo	
POS	z-sc.	POS	z-sc.
NOUN-DET-NOUN	-25.08	PRON-VERB-PRON	37.79
DET-NOUN-ADP	-23.62	PRON-VERB-ADV	33.92
NOUN-ADP-NOUN	-22.99	VERB-PRON-ADV	29.47
NOUN-ADP-DET	-22.96	VERB-PRON-PRON	21.79
DET-ADJ-NOUN	-22.71	PRON-PRON-VERB	19.78
ADP-DET-NOUN	-21.67	VERB-ADV-ADV	19.54
DET-NOUN-DET	-19.07	VERB-PRON-PRT	19.39
ADJ-NOUN-VERB	-18.33	PRON-VERB-PRT	19.25
ADP-DET-ADJ	-17.88	VERB-ADV-PRT	18.03
ADJ-NOUN-ADP	-17.64	PRON-VERB-ADJ	17.21

1. Preference for VERBS and PRON
2. Three very frequent POS triplets
3. Longer tail of infrequent POS
4. Steeper curve: not many mid-frequency POS

Frequency Distributions

