

Corpora are distributional data

- Corpus linguistics is a distributional discipline: everything is about
 - the **frequency of occurrence of something**
 - the frequency of **co-occurrence of something with something else**
 - both of the above **somewhere**
 - thus, we often report
 - **frequencies**
 - association measures (AMs)
 - dispersion measures (DMs, way too rarely!)
- what are those statistics used for, or grounded in?
- what motivates the use of such statistics?

Stefan Th. Gries
 UC Santa Barbara & JLU Giessen
<https://www.stgries.info>

The distributional hypothesis

- Central to much of corpus linguistics: different versions of the **distributional hypothesis**
- "you shall know a word by the company it keeps" (Firth (1957:11))
- "If we consider words or morphemes A and B to be more different in meaning than A and C, then we will often find that the distributions of A and B are more different than the distributions of A and C. In other words, **difference of meaning correlates with difference of distribution.**" (Harris 1970:585f.)
- see also Goldberg's (1995:50) semantic coherence principle, according to which constructions attract lexical items that are compatible with the semantic specifications of particular slots

In the context of Construction Grammar ...

- ... **how** is this typically done?
- **on the basis of 2x2-tables**
- for simple collexeme analysis

	construction	other	Sum
word _i	A	b	A+B
other	C	d	C+D
Sum	A+C	b+d	A+B+C+D

- for distinctive collexeme analysis

	construction ₁	construction ₂	Sum
word _i	A	B	a+b
other	C	d	C+D
Sum	A+C	B+D	a+b+C+D

- quantify the association of rows & columns to each other with typically, Φ -fisher-Yates exact test or G^2 (signed)
- sort word_i..._n by the AM

Corpora are distributional data

- Corpus linguistics is a distributional discipline: everything is about
 - the **frequency of occurrence of something**
 - the frequency of **co-occurrence of something with something else**
 - both of the above **somewhere**
 - thus, we often report
 - **frequencies**
 - association measures (AMs)
 - dispersion measures (DMs, way too rarely!)
- what are those statistics used for, or grounded in?
- what motivates the use of such statistics?

In the context of Construction Grammar ...

- ... **what** is typically done?
- the two by far most frequent applications are
 - **simple collexeme analysis**
 - probably most often, the association of a verb to a slot in an argument structure construction
 - **distinctive collexeme analysis**
 - the association of 'words' to one of several functionally similar constructions
 - probably most often, the association of a verb to a slot in two alternating constructions
- note: there's a reason for the quotes around *words*

Very wide applicability ...

- **These methods have been applied in a very wide variety of contexts**
 - general corpus-linguistic studies
 - psycholinguistic applications
 - second/foreign language acquisition/learning
 - diachronic studies of language change
 - experimental studies w/ colostr. results as predictors
 - experimental studies w/ colostr. results as controls
 - registers
 - languages (eg English, German, Dutch, Italian, Standard Arabic, Mandarin Chinese, ...)
- **and these methods have been hugely successful**
- but ...

Neglecting/ignoring dispersion

- Stefanwitsch & Gries (2003): verbs' attraction to the imperative in the ICE-GB: *let*, *see*, *look*, *listen*, *worry*, *verb₁*, *remember*, *check*, *verb₂*, *try*, *hang on*, ... makes perfect sense, but what are *verb₁* and *verb₂*?
- verb₂*: *fold*;
- 6th strongest with 16 occurrences in the imperative
- fold*->imperative: $KLD=2.095$ & *fold*-imperative: $KLD=0.031$
- log₁₀ P_{VE} of 21.03
- verb₂*: *process*;
- 9th strongest with 15 occurrences in the imperative
- proc*->imperative: $KLD=3.301$ & *proc*-imperative: $KLD=0.073$
- log₁₀ P_{VE} of 16.68
- in how many files of the ICE-GB are these found? $\frac{1}{500}$
- $KLD_{\text{diff}}^{\text{def}} for fold=0.9$; $KLD_{\text{diff}}^{\text{def}} for process=0.91$
- Gries (2024) on ditransitives: *tefl* is
- about as frequent in the ditransitive as *send*
- more strongly attracted to the ditransitive than *send*
- less well dispersed with the ditransitive than *send*

No more collocations, please!
 (at least not the traditional type →)

14

No more collocations, please!
 (at least not the traditional type →)

15

Collocations are monofactorial

- But it is linguistically/theoretically obvious that constructions involve more information than that: Fillmore (1988:36f.) distinguishes between internal and external properties of constructions:
- internal: formal, semantic, and symbolic structures
- external: recurring contextual patterns
- thus, structural knowledge of constructions is not the only kind of knowledge in speakers. Knowledge of contexts of use is also bona fide constructional knowledge.
- In a usage-based perspective, external properties range from co-textual ones (including colligational & collocational) to situational, social, and intertextual ones
- early example: Gries
- Haampe, & Schönefeld (2005): verbs are attracted to the *as*-predicative, but that attraction is in turn strongly attracted to the passive construction

No more collocations, please!
 (at least not the traditional type →)

16

No more collocations, please!
 (at least not the traditional type →)

17

Many collocational studies actually undermine their CXG background

- They do usually *not* quantify the attraction of a 'more grammatical construction' and a lexical construction ("it's constructions all the way down")
- they quantify the attraction of a 'more grammatical construction' with letter sequences! – why?
- because they treat a sequence of letters as a form side of a construction but not also its function
- Let me ask you this:
 - in a distinctive collexeme analysis of the dative alternation – ditransitives (*He gave him the book*) vs. prepositional datives (*He gave the book to him*) – *play* is attracted to the prepositional dative – how/why?
 - in a distinctive collexeme analysis of the verb-particle alternation – particles before DO (*He picked up the book*) vs. particle after DO (*He picked the book up*) – *carry out* prefers particle before DO
 - pick up* and *put down* have no strong preference – how/why?

No more collocations, please!
 (at least not the traditional type →)

18

No more collocations, please!
 (at least not the traditional type →)

19

Collocations are monofactorial, given that ...

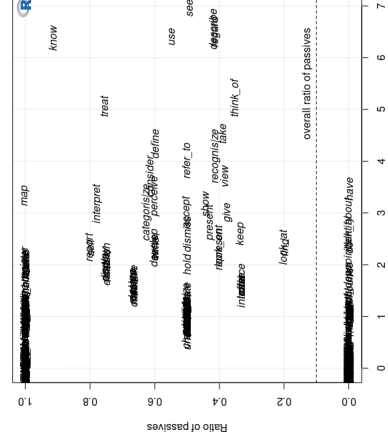
- ... they are computed from 2x2 tables cross-tabulating a construction as a function of a lexical predictor
- ... they can be re-written as a logistic regression: $g\text{lm}er(\text{CONSTRUCTION} \sim 1 + (1|\text{WORD}))$
- but it is statistically obvious that a constructional choice will not be due only to lexical items
- but it is linguistically/theoretically obvious that constructions involve more information than that:
 - On a usage-based model then, constructions are simply conventionalized chunks of linguistic knowledge... From this, it follows that the storage and organization of grammatical knowledge is dependent on, and can change according to, patterns of use" (Patten 2014)

No more collocations, please!
 (at least not the traditional type →)

Stefan Th. Gries
 UC Santa Barbara & JLU Giessen

15

Collocations are monofactorial



Logged frequency of verb lemma everywhere

No more collocations, please!
 (at least not the traditional type →)

Stefan Th. Gries
 UC Santa Barbara & JLU Giessen

17

How certain, robust, variable are the results?

- There is a notable overall increase of statistical sophistication in linguistics
- more significance testing, more effect sizes, more confidence intervals, more Bayesian methods, ...
- but collocational applications are lagging behind
- some measure 'are' (essentially) p -values, ok, but they are usually not corrected for post-hoc testing (Gries 2005, CLLT)
- they arguably run into problems with the random-sampling assumption of null hypothesis significance testing
- some measures are effect sizes, but
- we do not get (equivalents of) confidence intervals, which obfuscates how much trust we can place in the ranked lists we usually interpret

No more collocations, please!
 (at least not the traditional type →)

Stefan Th. Gries
 UC Santa Barbara & JLU Giessen

19

The reviewers: revise & resubmit 5: constructions all the way? on meaning & polysemy
 The editor: Here's my ranking of the suggestions 2: please! 3: please! 4: please!
 Concluding remarks: given a traditional analysis, a reviewer should ...

Introduction 4: essentially monofactorial/flat view of usage
 The reviewers: revise & resubmit 2: conflating both directions of attraction
 The editor: Here's my ranking of the suggestions 2: please! 3: please! 4: please!
 Concluding remarks: essentially monofactorial/flat view of usage

If you really submit a traditional collexeme analysis in 2026, ...

- ... you should get maximally a *revise and resubmit* ;-)
- reviewers should say
 - tupletize
 - measure frequency
 - measure association without frequency
 - measure association directionally
 - measure dispersion without frequency
 - add other variables
 - take constructions and their senses more seriously
 - add uncertainty assessments
- Let's look at what this can add to your work

No more colstructions, please!
 (at least not the traditional type ->)

Stefan Th. Gries
 UC Santa Barbara & JLU Giessen

No more colstructions, please!
 (at least not the traditional type ->)

Stefan Th. Gries
 UC Santa Barbara & JLU Giessen

Introduction 1: conflating frequency & association
 The reviewers: revise & resubmit 2: conflating both directions of attraction
 The editor: Here's my ranking of the suggestions 2: please! 3: please! 4: please!
 Concluding remarks: essentially monofactorial/flat view of usage

Introduction 1: conflating frequency & association
 The reviewers: revise & resubmit 2: conflating both directions of attraction
 The editor: Here's my ranking of the suggestions 2: please! 3: please! 4: please!
 Concluding remarks: essentially monofactorial/flat view of usage

why one needs to control for frequency

	LOW	Cof-ditrans	Cof-ther	Sum	
V-assured	0.1	12.1	12.2		
Weather	182.0	137075.1	138895.2		
Sum	182.0	137087.2	138907.4	-0.6769	-3.2618
OIS					
Weather	182.0	137085	138895		
Sum	182.0	137087	138907	5.9910	4.2341
UPP					
Weather	182.0	137087	138907		
Sum	182.0	137087	138907	6.2540	6.2540

- PMI theoretically falls into the interval $[-\infty, \infty]$**
- but practically, ie given marginal totals**
 - PMI for assured & ditrans falls into $[-0.68, 6.25]$
 - PMI for offer & ditrans falls into $[-3.26, 6.25]$
- seeing this, can we really compare
- assured's PMI of 5.99 to
- offer's PMI of 4.23?
- Look what can actually happen ...

No more colstructions, please!
 (at least not the traditional type ->)

Stefan Th. Gries
 UC Santa Barbara & JLU Giessen

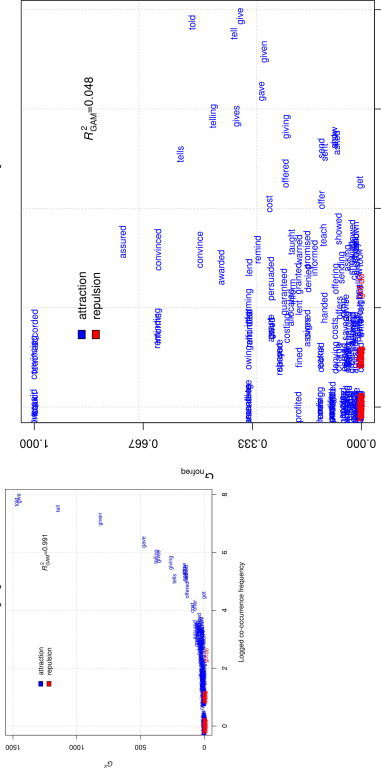
No more colstructions, please!
 (at least not the traditional type ->)

Stefan Th. Gries
 UC Santa Barbara & JLU Giessen

Introduction 1: conflating frequency & association
 The reviewers: revise & resubmit 2: conflating both directions of attraction
 The editor: Here's my ranking of the suggestions 2: neglect dispersion
 Concluding remarks: essentially monofactorial/flat view of usage

Introduction 1: conflating frequency & association
 The reviewers: revise & resubmit 2: conflating both directions of attraction
 The editor: Here's my ranking of the suggestions 2: neglect dispersion
 Concluding remarks: essentially monofactorial/flat view of usage

what happens if one 'takes freq out' ...



No more colstructions, please!
 (at least not the traditional type ->)

Stefan Th. Gries
 UC Santa Barbara & JLU Giessen

No more colstructions, please!
 (at least not the traditional type ->)

Stefan Th. Gries
 UC Santa Barbara & JLU Giessen

"The authors should tupletize to keep dimensions of information separate"

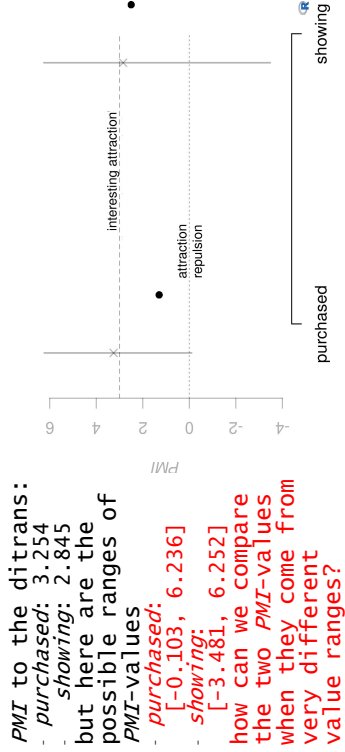
- "They should not compute one measure that conflates a lot of frequency and association (p_{VE} , G^2 , chi-squared, t , Z , ...)
- one measure that is 'only' (similar to) an effect size (PMI , (log) OR, ...)
- they should determine/compute frequencies of (co-)occurrence
- one measure of association for each direction, but make that measure
 - one that distinguishes attraction from repulsion
 - one works the same in each direction
 - one that is not/little affected by frequency
- In one case study, the most widely-used AM - G^2 - is $R^2_{GAM}=0.9465$ correlated with co-occurrence frequency
- In another (see above), G^2 is $R^2_{GAM}=0.991$ correlated with co-occurrence frequency
- the authors might use Gries's (2024) suggestion and compute the KL-divergence in both directions
- use the version that 'controls' for frequency

No more colstructions, please!
 (at least not the traditional type ->)

Stefan Th. Gries
 UC Santa Barbara & JLU Giessen

Introduction 1: conflating frequency & association
 The reviewers: revise & resubmit 2: conflating both directions of attraction
 The editor: Here's my ranking of the suggestions 2: please! 3: please! 4: please!
 Concluding remarks: essentially monofactorial/flat view of usage

why one needs to control for frequency



PMI to the ditrans:
 purchased: 3.254
 showing: 2.845
 but here are the possible ranges of PMI-values
 purchased: [-0.103, 6.236]
 showing: [-3.481, 6.252]
 how can we compare the two PMI-values when they come from very different value ranges?

No more colstructions, please!
 (at least not the traditional type ->)

Stefan Th. Gries
 UC Santa Barbara & JLU Giessen

Introduction 1: conflating frequency & association
 The reviewers: revise & resubmit 2: conflating both directions of attraction
 The editor: Here's my ranking of the suggestions 2: neglect dispersion
 Concluding remarks: essentially monofactorial/flat view of usage

"The authors should tupletize to keep dimensions of information separate"

- They should check for all co-occurrences the degree to which they are widely dispersed
- however, here, they must be even more aware of the high correlation of nearly all DMS with frequency
- this is esp. true of range, but also other DMS
- the authors might use Gries's (2024) suggestion and compute the KL-divergence
 - of the distribution of co-occurrences (posterior) from the file sizes (prior)
- that controls for frequency
- that is the measure that can reflect dispersion differences between lower-frequency words
- has been shown to boost predictive modeling of reaction times best (Gries 2022, JSL5)

No more colstructions, please!
 (at least not the traditional type ->)

Stefan Th. Gries
 UC Santa Barbara & JLU Giessen

No more colstructions, please!
 (at least not the traditional type ->)

Stefan Th. Gries
 UC Santa Barbara & JLU Giessen

"The authors should consider more than just letter sequences"

- Above, I said most studies treat a sequence of letters as a construction ...
- ... when a sequence of letters is only the form side of a construction but not also its function
- then, I asked why is/are
- *play out* attracted to the prepositional dative
- *carry out* attracted to verb-particle-DO (VPO)?
- *pick up* and *put down* not attracted to either construction?
- you see what this presupposes?
- to answer some of these, we need to answer:
 - what is the construction [*carry out*]?
 - what is the construction [*pick up*]?
 - what is the construction [*put down*]?

Constructions, not letter sequences: [verb particle] + 'meaning'!

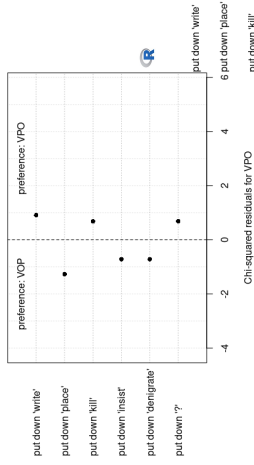
- [*carry out*] is the formal pole – but the functional?
 - 'move outside', (eg *a box out of an apartment*)
 - 'execute/perform', (eg *an order*)
- [*pick up*] is the formal pole – but the functional?
 - 'register/notice/learn', (eg *I don't pick things up quite as quickly as he does*)
 - 'take', (eg *He picked up a pencil*)
 - 'meet/hit on', (eg *He picked a girl up on the train*)
 - 'address', (eg *Let's pick up the point about membership fees again*) ...
- [*put down*] is the formal pole – but the functional?
 - 'kill', (eg *put down the terminally ill dog*)
 - 'place', (eg *put down the cup*)
 - 'write/register', (eg *put me down for some pudding*)
 - 'denigrate', (eg *must you put him down like that? He's only just started practising this!*)

Collocations require a different annotation

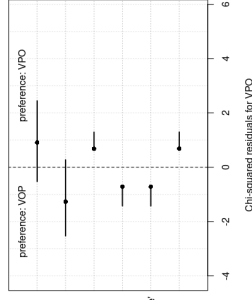
- "Here's what everyone is doing collocations on
- "here's what we should be doing collocations on
- "what are the results?
 - for the 2 constructions of the letter sequence *carry out*
 - [*carry out*] 'execute': 49 VPO : 0 VOP
 - [*carry out*] 'move out': 0 VPO : 1 VOP
 - for the 6 constructions of the letter sequence *put down*
 - it is **self-contradictory** to
 - self-identify as **construction** Grammarians
 - do **collocations**
 - take seriously only the **constructionhood** – a pairing of **form and function** – of one of the constructions in one's **grammatical**/argument structure construction
 - the ditransitive or prepositional dative
 - verb-particle-DO or verb-DO-particle
 - **but reduce the other one to the form side** of 'what letters are being used' and disregard their meaning!

Collocations require a different annotation

Distinctive put down-collexemes of VPO



Distinctive put down-collexemes of VPO



	A	B	C	D	E	F	G	H
1	PRECEDING	VERB	DO	PART	SUBSEQUENT	CX	LEMMA-LETTER	
2	Mid of a book and you just don't want to	put	it	down	on my C.V. though	VOP	put down	
3	God you really know how to	put	someone	down	don't you	VOP	put down	
4		put	my foot	down	at being a medic, so this is a compromise	VOP	put down	
5		put	them	down	there	VOP	put down	
6	You said well fuck off then -laughter- and	put	the phone	down	on him you know	VOP	put down	
7	Now on my uh my manifesto do I	put	* create and let	down	cos when I'm when I'm listing my secop	VPO	put down	
8	To move your hand and then I	put	it	down	by your side cos you've only half used	VOP	put down	
9	Well I puting	put	it	down	by your side cos you've only half used	VOP	put down	
10	and then you	put	it	down	by your side cos you've only half used	VOP	put down	
11	To move your hand and then I	put	it	down	by your side cos you've only half used	VOP	put down	
12	Uhrr now I	put	a lie	down	the whole point is which one is right	VOP	put down	
13	I'm an assistant manager and they	put	both these	down	and just present it so that it's all pages	VPO	put down	
14	So you	put	dived	down	on the form with the new address	VPO	put down	
15	So you	put	* names as the	down	on the form with the new address	VPO	put down	
16	All you can really do then is just is just	put	* but it's not be	down	alm	VPO	put down	
17	And just	put	the cost prices	down		VPO	put down	
18								
19								
		VPO	VOP					
		put down	a	b				
		other	c	d				
		Sum	a+b	c+d				
			a+b+c+d					

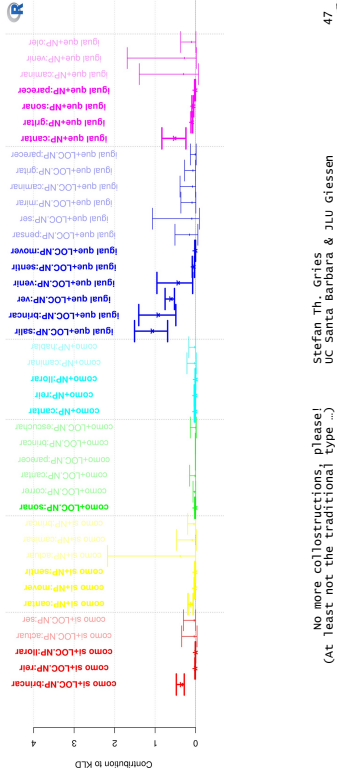
Some collocation studies did much better

- **wiechmann (2008)** found that collocation results that distinguished senses of verbs (whether they take a nominal or a sentential complement) outperform collocation results that did not distinguish senses when it comes to predicting reading latencies
- **Gilquin (2013)** showed different senses of verbs have v different preferences for causative constructions
- **Colleman & Bernolet (2012)** showed how polysemy effects can partially explain the lack of (better) agreement between corpus & exp. data on the dative alternation
- **Bernolet & Colleman (2016)** illustrated how verb senses can have widely different alternation biases (in a study of the Dutch dative alternation) and how, in an experiment, priming was affected by an interaction of the strength of priming and sense-specific biases of target verbs

The editor: Here's my ranking of suggestions in terms of importance
 The reviewers: revise & resubmit
 Concluding remarks

Introduction ranking in terms of simplicity
 The editor: Here's my ranking of suggestions in terms of importance
 The reviewers: revise & resubmit
 Concluding remarks

"As an editor, ..."



No more collocations, please!
 (At least not the traditional type ...)
 Stefan Th. Gries
 UC Santa Barbara & JLU Giessen

Logged co-occurrence freq.
 (At least not the traditional type ...)
 Stefan Th. Gries
 UC Santa Barbara & JLU Giessen



The editor: Here's my ranking of suggestions in terms of importance
 The reviewers: revise & resubmit
 Concluding remarks

Introduction ranking in terms of simplicity
 The editor: Here's my ranking of suggestions in terms of importance
 The reviewers: revise & resubmit
 Concluding remarks

"As an editor, ..."

- However, it is also necessary to consider the importance of the suggestions in terms of how much they are required and improve the analysis
- For that, the ranking is probably a bit different:
 - **annotate senses**
 - because otherwise you're not exploring co-occurrence of constructions
 - **distinguish frequency and association and different directions of association**
 - because otherwise you're not 'measuring' what you say you're 'measuring'
 - **annotate other variables**
 - because of course there's more to the uses of these constructions than just the words/lemmas in one slot
 - **compute dispersion**
 - because you need to at least make sure you're not overinterpreting outlier results
 - **compute ranges of plausible values**
 - because it's just one loop on top of dispersion and will help streamlining the interpretation

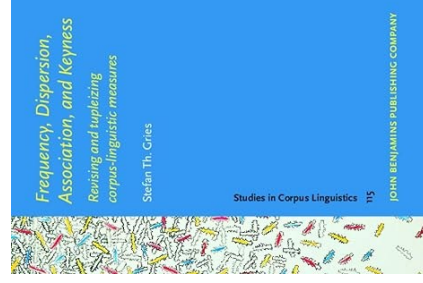
- Collocations is a general insightful approach, ...
- but linguistics as a discipline has evolved a lot when it comes to the quantitative/statistical analysis of (corpus) data
- it's not 2003/2004 anymore
- just like variationist sociolinguistics moved on from varbrul (at some point, finally, thankfully ...), ...
- ... constructional studies, and quantitative approaches to CXG in general, need to move on from 'basic collocations'
- collocations need to be 'updated' and I hope this talk has shown you why and especially how
- so, go and **preempt this reviewer!** ;-)

Take-home message

No more collocations, please!
 (At least not the traditional type ...)
 Stefan Th. Gries
 UC Santa Barbara & JLU Giessen

No more collocations, please!
 (At least not the traditional type ...)
 Stefan Th. Gries
 UC Santa Barbara & JLU Giessen

Introduction ranking in terms of simplicity
 The editor: Here's my ranking of suggestions in terms of importance
 The reviewers: revise & resubmit
 Concluding remarks



Thank you!

<https://www.stgries.info>