

On frequency in corpora 1: frequencies vs. association measures

Stefan Th. Gries
Department of Linguistics
University of California, Santa Barbara
<http://tinyurl.com/stgries>

Collocation analysis: its goals, computation, results, and interpretation

- Collocation analysis (CA) applies association measures (AMs) to lexico-syntactic co-occurrence
 - 1: retrieve all instances of a construction c
 - 2: compute AMs for all collexemes w in a slot of c
 - 3: rank collexemes w by their association score
 - 4: explore the top t collexemes' functional patterns

	w	$\neg w$	totals
c	a	b	a+b
$\neg c$	c	d	c+d
totals	a+b	b+d	a+b+c+d

- any association measure can be used (S&G 2003:217)
 - by far most frequent approach: $(-)\log_{10} p_{\text{Fisher-Yates exact}}$
 - no distributional assumptions
 - can handle the low-frequency data following from Zipf
 - as a p -value, p_{FYE} is correlated with both effect size and sample size, weighing an observed effect higher, when it is found in a larger sample: ' $p_{\text{FYE}}^{14/35} > p_{\text{FYE}}^{8/20}$ '

The logic of collocations: (distinctive) collexeme analysis :-|

- Collocations are but one approach towards co-occurrence data, which is based on a statistical test of bi-directional association in 2x2 tables
 - **collexeme analysis** quantifies the degree of attraction/repulsion of one construction (a word) to another
 - **distinctive collexeme analysis** quantifies the degree of attraction/repulsion of one construction (a word) to one of x functionally similar constructions
- crucially, the approach
 - serves to **rank-order collexemes**: the top collexemes reflect the construction's prototypical sense(s)
 - **normalizes frequency of occurrence** in a construction
 - typically uses $-\log_{10} p_{\text{Fisher-Yates exact}}$ – because it reflects **association, significance, and frequency** (cf. below)
 - can distinguish **attracted and repelled collexemes**
 - CA has been applied to data from English, German, Dutch, Swedish, ... synchronic and diachronic data, syntactic alternations, priming, second language acquisition, ...

"Problems of collostructional analysis" (Bybee 2010: Section 5.12)

- Bybee criticizes the collostructional approach for what she considers "problems"
 - (1) the fact that a **significance test** is used because "lexemes do not occur in corpora by pure chance" and "the factors that make a lexeme high frequency in a corpus [may be] the factors that make it a central and defining member of the category" (p. 97) And, how is d ($\neg w$, $\neg c$) calculated? (p. 98)
 - (2) no **"cognitive mechanism"** is proposed that "corresponds to their analysis" (p. 100)
 - (3) "Since no semantic considerations go into the analysis, it seems plausible that **no semantic analysis can emerge from it**" (p. 98) ... "since it works only with numbers and not with meaning." (p. 100)
 - (4) collostructions do not distinguish **low-frequency semantically related and semantically unrelated collexemes** in the data of Bybee & Eddington (2006)

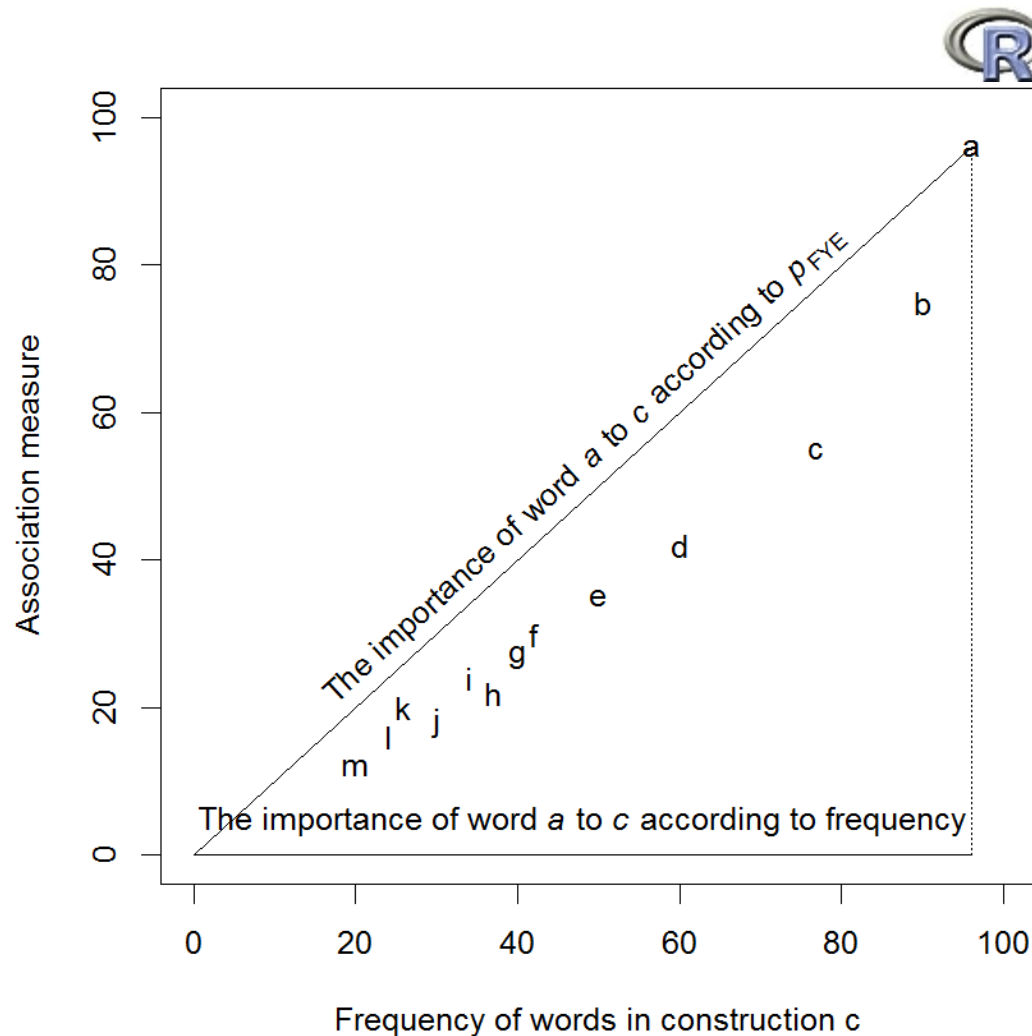
- 1: using a significance test and non-robustness Any association measure may work and will normalize
- 2: cognitive mechanisms Some advantages of Fisher-Yates exact tests
- 3: collocations and semantic analysis Nonlinearity and the robustness of rankings
- 4: the discriminative power of collocations Other studies, other measures

On (1) and what association measures do

- (1) misses the point completely
 - no corpus linguist would ever say words occur by chance
 - if they did, not only computing AMs but also counting frequencies would be pointless ...
 - the use of a significance test is merely one statistical heuristic:
 - S&G (2003:217) explicitly said *any AM* could be used
 - measures not based on *p*-values have in fact been used
 - what this approach - any AM - does is
 - downgrade words that are highly promiscuous
 - upgrade words that are highly faithful to the construction under investigation
 - thus
 - it makes *regard* - not *see*, *describe*, or *know* - most representative of the *as*-predicative (V NP_{DO} *as* complement)
 - it recognizes that equal frequencies of occurrence may mask preferences (*consider* is equally frequent in *to*- and *ing*-complement constructions, but *to*- is much more freq)

- 1: using a significance test and non-robustness Any association measure may work and will normalize
- 2: cognitive mechanisms Some advantages of Fisher-Yates exact tests
- 3: collocations and semantic analysis Nonlinearity and the robustness of rankings
- 4: the discriminative power of collocations Other studies, other measures

On (1) and what association measures do

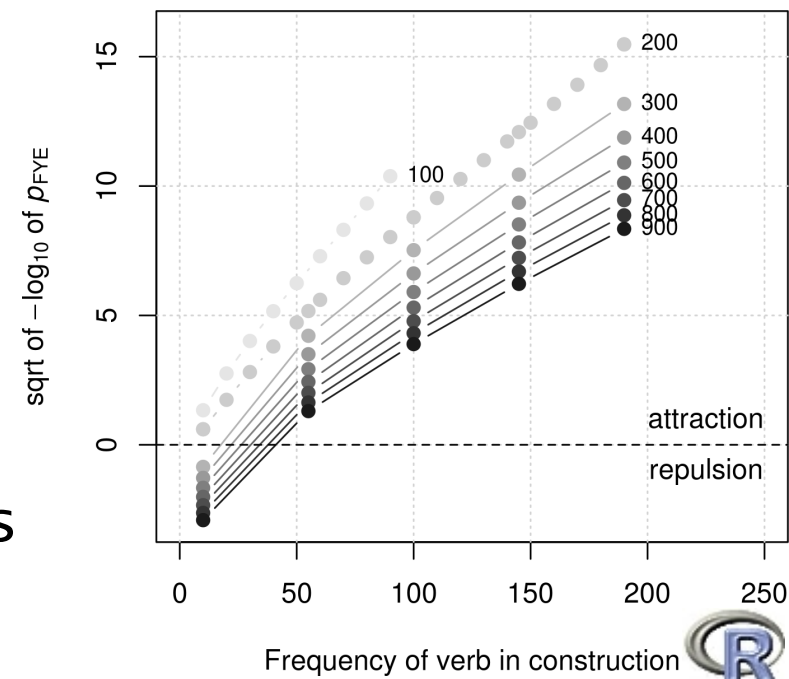


- 1: using a significance test and non-robustness Any association measure may work and will normalize
- 2: cognitive mechanisms Some advantages of Fisher-Yates exact tests
- 3: collocations and semantic analysis Nonlinearity and the robustness of rankings
- 4: the discriminative power of collocations Other studies, other measures

The most wide-spread measure of collocation strength

The Fisher-Yates exact test (FYE) is

- statistically more appropriate than many other measures which make unwarranted distributional assumptions
- a significance test but unlike many other measures
 - it actually incorporates frequencies to make the measure more compatible with frequency/entrenchment-based accounts (cf. Stefanowitsch & Gries 2003)
 - with increasing frequency of co-occurrence, the collexeme strength goes up more (both is what Bybee would want!)
- statistics need not be simple to provide accurate representations of distributional properties of language (Stefanowitsch 2012)
- Wiechmann's (2008) comparison ranks FYE second (after *MinSem*, which is theoretically problematic)



- 1: using a significance test and non-robustness Any association measure may work and will normalize
- 2: cognitive mechanisms Some advantages of Fisher-Yates exact tests
- 3: collocations and semantic analysis Nonlinearity and the robustness of rankings
- 4: the discriminative power of collocations Other studies, other measures

which ranking of ditransitive verbs do you prefer?

Verbs	Fisher-Yates exact	Verbs	Log odds	Verbs	Mutual Information
<i>give</i>	308.25	<i>give</i>	2.20	<i>accord</i>	6.07
<i>tell</i>	126.80	<i>accord</i>	2.13	<i>award</i>	5.87
<i>send</i>	67.14	<i>award</i>	2.02	<i>give</i>	5.73
<i>offer</i>	48.48	<i>allocate</i>	1.73	<i>allocate</i>	5.26
<i>show</i>	32.65	<i>profit</i>	1.65	<i>profit</i>	5.07
<i>cost</i>	21.95	<i>owe</i>	1.63	<i>owe</i>	5.01
<i>teach</i>	15.36	<i>lend</i>	1.59	<i>lend</i>	4.92
<i>award</i>	10.86	<i>offer</i>	1.58	<i>offer</i>	4.86
<i>allow</i>	9.95	<i>cost</i>	1.52	<i>cost</i>	4.72
<i>lend</i>	8.55	<i>send</i>	1.51	<i>grant</i>	4.69
<i>deny</i>	8.35	<i>grant</i>	1.50	<i>send</i>	4.63

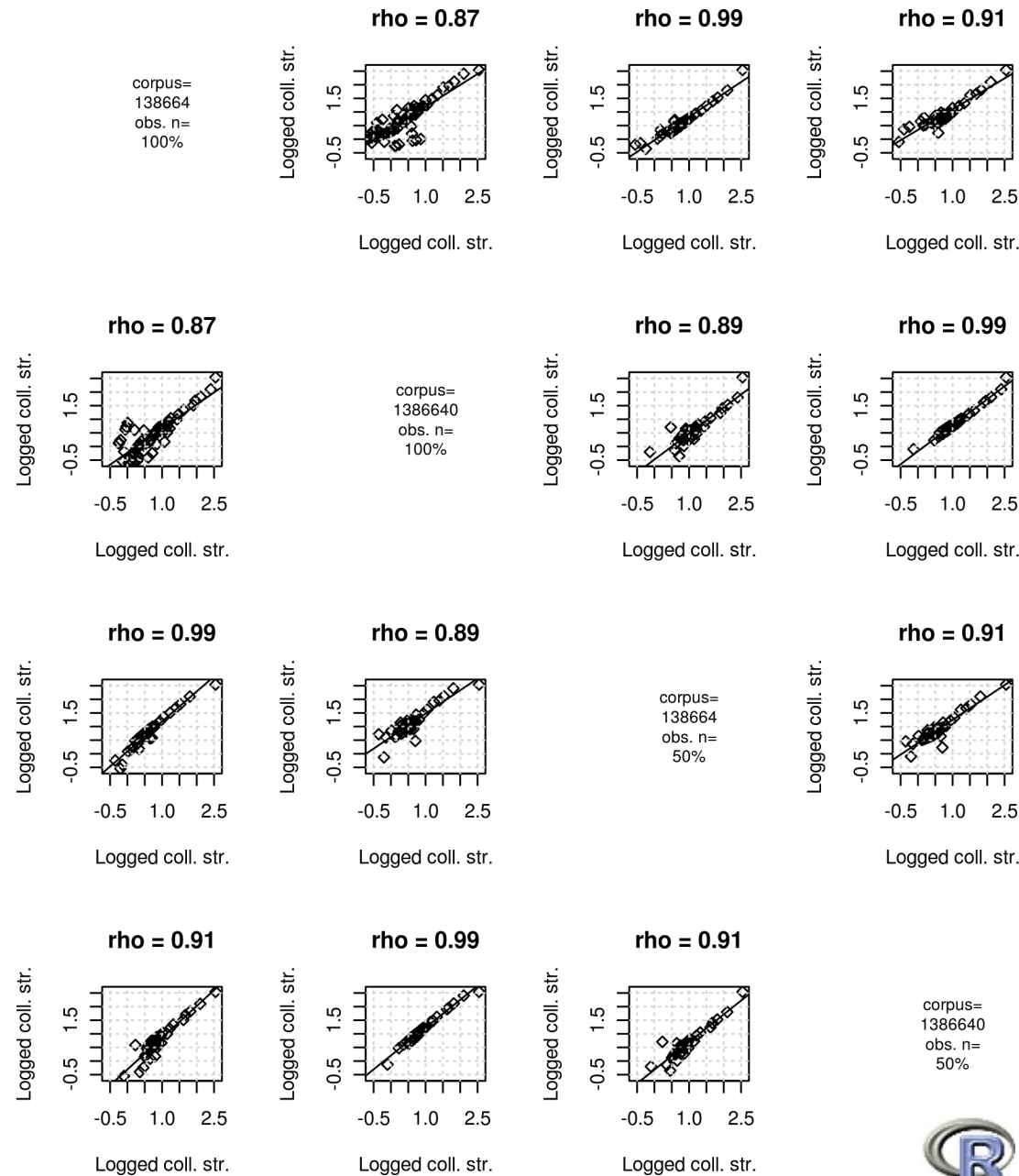
- 1: using a significance test and non-robustness Any association measure may work and will normalize
- 2: cognitive mechanisms Some advantages of Fisher-Yates exact tests
- 3: collocations and semantic analysis Nonlinearity and the robustness of rankings
- 4: the discriminative power of collocations Other studies, other measures

On (1) and what association measures do

- (1) misses the point completely
 - association measures actually have a variety of characteristics that should appeal to cognitive linguists
 - association measures typically result in frequency-and-their-effects relations that are **non-linear**, just as frequency effects (Tryk 1968), learning (Anderson 1982), forgetting curves, priming decay (Gries 2005, Szmrecsanyi 2005), etc.
 - association measures reflect
 - that "[r]aw frequency of occurrence is less important than the contingency between cue and interpretation", and
 - the notion that "Contingency, and its associated aspects of predictive value, information gain, and statistical association, have been at the core of learning theory ever since." (Ellis & Ferreira-Junior 2009)
 - d is computed as it is for all AMs: on the basis of a reasonable approximation for a corpus; simulations show that CA rankings are very robust

- 1: using a significance test and non-robustness Any association measure will work and will normalize
- 2: cognitive mechanisms Some advantages of Fisher-Yates exact tests
- 3: collocations and semantic analysis Nonlinearity and the robustness of rankings
- 4: the discriminative power of collocations Other studies, other measures

On (1) and what AMs do



On frequency in corpora 1:
frequencies vs. association measures

Stefan Th. Gries
University of California, Santa Barbara



- 1: using a significance test and non-robustness Any association measure may work and will normalize
- 2: cognitive mechanisms Some advantages of Fisher-Yates exact tests
- 3: collocations and semantic analysis Nonlinearity and the robustness of rankings
- 4: the discriminative power of collocations Other studies, other measures

So, what if we pit frequency against association strength?

- Gries, Hampe, & Schönefeld (2005): *as*-predicative
 - *regard* (80) is returned as the central verb (see (111) and *know* (79) are more/as frequent but more general than the top-ranked *regard* and *describe*)
 - in a sentence completion task, p_{FYE} outperforms conditional prob. and freq. as predictors of completions
- Ellis & Simpson-Vlach (2009): *MI* outperforms frequency as a predictor of formulaicity ratings
- Gries, Hampe, & Schönefeld (2010): in a self-paced reading times experiment, p_{FYE} ($p=0.065$) outperforms freq ($p=0.293$) with an effect size 3 times as high
- Coleman & Bernolet (2012) find seemingly erratic verb-specific preferences in the Dutch dative alternation that fall into place when explored with association strength

- 1: using a significance test and non-robustness Any association measure may work and will normalize
- 2: cognitive mechanisms Some advantages of Fisher-Yates exact tests
- 3: collocations and semantic analysis Nonlinearity and the robustness of rankings
- 4: the discriminative power of collocations Other studies, other measures

$-\log_{10} p_{\text{FYE}}$ yields good results,
but can we do better than that?

frequency of (!)x and (!)y in some corpus	2		Totals
	y	!y	
1: x	a	b	a+b
1: !x	c	d	c+d

- Nearly all AMs are bidirectional – but (associative) learning is not
- thus, a **uni-directional AM** may be more useful

$$\Delta P_{2|1} = p(\text{word}_2 | \text{word}_1 = \text{present}) - p(\text{word}_2 | \text{word}_1 = \text{absent}) = \frac{a}{a+b} - \frac{c}{c+d}$$

$$\Delta P_{1|2} = p(\text{word}_1 | \text{word}_2 = \text{present}) - p(\text{word}_1 | \text{word}_2 = \text{absent}) = \frac{a}{a+c} - \frac{b}{b+d}$$

- consider *of course* in the spoken part of the BNC
 - $MI=5.41$, $t=476.97$, $G^2=36,693.85$, $p_{\text{FYE}} < 10^{-320}$
 - $\Delta P(\text{course} | \text{of}) = 0.032$ and $\Delta P(\text{of} | \text{course}) = 0.697$
- discrepancies like this are rather common
- a reanalysis of Gries, Hampe, & Schönefeld (2005) with ΔP (construction|verb) shows that ΔP is a significant predictor of subjects' completions

- 1: using a significance test and non-robustness
- 2: cognitive mechanisms
- 3: collocations and semantic analysis
- 4: the discriminative power of collocations

On (2), cognitive mechanisms that may be involved / have been invoked

- In addition to the above comments, collocational studies by Stefanowitsch and myself as well as others and other work by myself at least have embraced particular cognitive mechanisms
 - Stefanowitsch & Gries (2003:237)
 - discuss implications of collocations to [psycholinguistic studies of language acquisition](#) (quoting Goldberg 1999)
 - discuss the relation of collocation strength to raw token frequency in that connection
 - draw an explicit connection to the notion of [cue validity](#), which is a well-established notion in the Competition Model, Goldberg's work, Ellis's work, ...
 - Stefanowitsch & Gries (2003:239, n. 6) discuss the relation of collocation strength to [entrenchment](#) in a way that is in fact sympathetic to Bybee (cf. also Stefanowitsch 2008 for more discussion)
 - Wiechmann (2008:257) discusses association strengths (incl. collocations) and their relation to [cue validity and cue strength](#) (cf. also MacWhinney p.c.)

- 1: using a significance test and non-robustness
- 2: cognitive mechanisms
- 3: collocations and semantic analysis
- 4: the discriminative power of collocations

On (3) and on semantic analysis in general

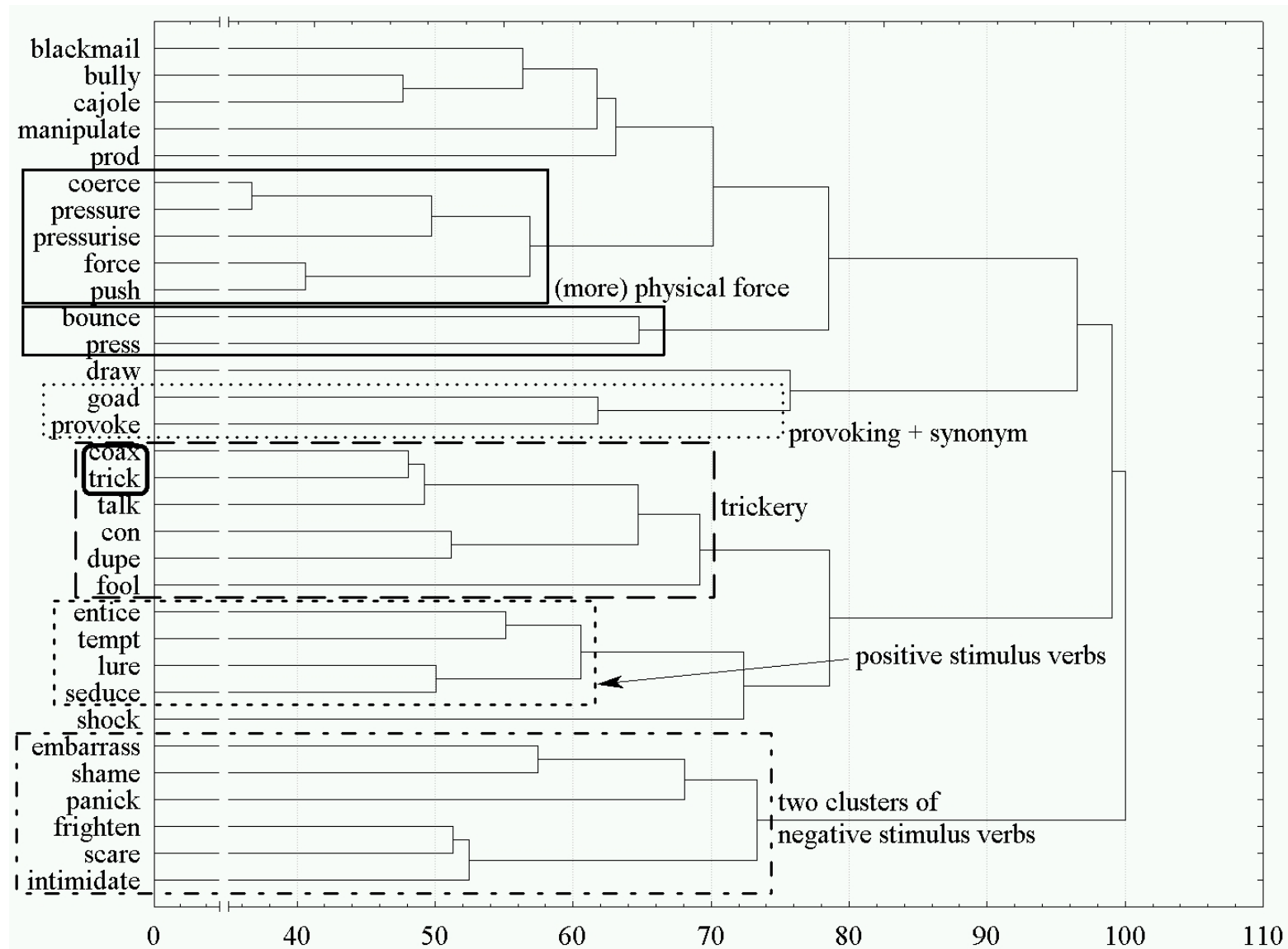
- In general, the notion that semantics cannot emerge from analysis X if no semantics is entered into it is flawed (and I know that Joan knows that)
 - there is a whole body of work in computational (psycho)linguistics, or distributional linguistics, where purely number-based **distributional analyses reveal functionally highly coherent clusters**
 - Reddington, Chater, & Finch (1998)
 - input: co-occurrence frequencies of 150 bigrams before 1000 target words
 - output: cluster analysis that returns parts of speech
 - Mintz, Newport, & Bever (2002): similar analysis but different definition of context (boundaries at function words)
 - a lot of work in information retrieval, document summarization, latent semantic analysis, etc.

On (3) and semantic analysis in collostructional analysis

- In particular, in collostructional analysis semantic analysis follows the statistics, but it is false to assume that no semantic analysis can emerge from it
 - there are many studies now on very many different constructions and 'alternations' that have shown this, e.g. Stefanowitsch & Gries's (2003) on the ditransitive
 - the verb returned as prototypical/central is *give*
 - the next verbs instantiate the senses of Goldberg's (1992, 1995) polysemy analysis of the ditransitive: satisfaction conditions (*offer, owe, promise*), enablement (*allow*), cause non-transfer (*deny*), cause future transfer (*grant*), intention to receive (*earn*), communication as transfer (*tell, teach*), perceiving as receiving (*show*), ... (cf. above)
 - Gries & Stefanowitsch's (2010) cluster
 - verbs in the *into*-causative based on the *ing*-verbs
 - verbs in the *way*-construction based on the prepositions

- 1: using a significance test and non-robustness Semantics as emerging from distributional linguistics
- 2: cognitive mechanisms Semantics as emerging from collostructional analysis
- 3: collostructions and semantic analysis Co-varying collexemes in the *into*-causative
- 4: the discriminative power of collostructions Co-varying collexemes in the *way*-construction

V_{action} NP_{Di}rob Patient *into* **V_{ing}**

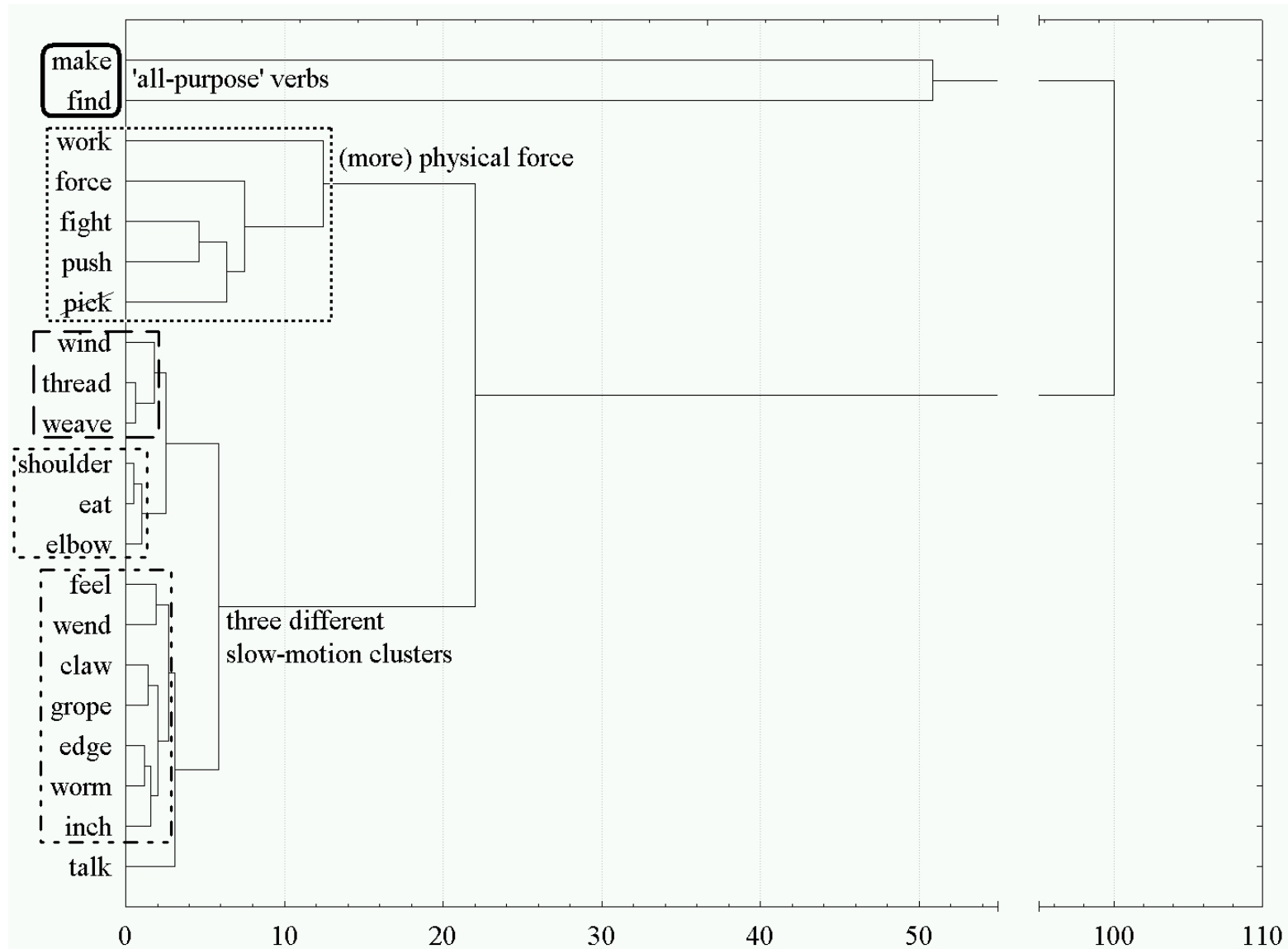


On frequency in corpora 1:
frequencies vs. association measures

Stefan Th. Gries
University of California, Santa Barbara

- 1: using a significance test and non-robustness
 - 2: cognitive mechanisms
 - 3: collocations and semantic analysis
 - 4: the discriminative power of collocations
- Semantics as emerging from collocational analysis
Co-varying collexemes in the *into*-causative
Co-varying collexemes in the *way*-construction

V POSS *way* [_{PP} **PREP** NP]



On frequency in corpora 1:
frequencies vs. association measures

Stefan Th. Gries
University of California, Santa Barbara

- 1: using a significance test and non-robustness
- 2: cognitive mechanisms
- 3: collocations and semantic analysis
- 4: the discriminative power of collocations

On (4) and (the perceived lack of) discriminative power

- Bybee: low-frequency collexemes do not distinguish between semantically related and semantically unrelated items (referring to Bybee & Eddington 2006)
- how does she 'show' that?

Step	CA	Bybee's 'CA'
1	retrieve all instances of <i>c</i>	-
2	compute AMS for all collexemes in a slot of <i>c</i>	compute AMS for 24 semantically-analyzed adjectives w/ freqs of ≤1
3	rank-order collexemes by the AMS	-
4	explore top <i>t</i> collexemes functionally	-

- by ignoring all previously published results
- by ignoring how collexeme analysis is actually done
 - semantic analysis is done **after** the stats – not before
 - **all** collexemes are included – not just those with $n \leq 1$
- note, she focuses only on whether collocations can predict her positive judgments – what about negative judgments? (Gries & Wulff 2009)

- 1: using a significance test and non-robustness what is it that Bybee calls *collostructional analysis*?
 2: cognitive mechanisms Results of a real collostructional analysis (recap)
 3: collostructions and semantic analysis
 4: the discriminative power of collostructions

On (4) and (the perceived lack of) discriminative power

- what do collostructions output when applied as intended – inspecting words ranked highest based on a measure combining frequency and contingency?
 - S & G (2003) and G & S (2004)
 - ditransitives, dative alternation, *as*-predicatives: as above
 - *into*-causatives: as above, and more comprehensive results than Hunston & Francis's (1999) frequency-based results
 - Ellis & Ferreira-Junior (2009): freq. of learner uptake is predicted by frequency, ΔP , $-\log_{10}(p_{\text{FYE}})$ – in fact, $-\log_{10}(p_{\text{FYE}})$ outperforms freq in $2/3$ constructions
 - Coleman & Bernolet (2012): verb-specific preferences (Dutch dative alternation) that appear erratic from a raw frequency perspective turn out to be systematic once distinctive collexemes are used
 - Gregory et al. (1999): MI is correlated significantly with 3 pronunciation effects (frequency with 2)

Some additional difficulties I have, and where that leads us

- Bybee criticizes collocations for ignoring **low-frequency collexemes**, which may reflect productivity
 - but, again, that's not what collocations are supposed to do: if one wants a productivity measure, one can count the number of hapaxes directly from the clx output (but note that only clx can rank low-freq words highly)
- Bybee argues that it is not pertinent how often a lexeme does *not* occur in a construction, but
 - we have seen the literature on learning says otherwise
 - CA just uses one more type of information than Bybee herself (9)

		$\left\{ \begin{array}{l} me \\ you \\ him \\ her \\ the producer \end{array} \right\}$	$\left\{ \begin{array}{l} mad \\ crazy \\ up the wall \end{array} \right\}$
SUBJECT	[DRIVE]		

CA just normalizes the frequency of *crazy/mad/etc.* against their overall frequencies

Some additional difficulties I have, and where that leads us

- Her example (9) simplifies but is instructive!
- ultimately, the situation is more complicated than Bybee AND collocations assume: we really need
 - the **type frequencies** and **token frequencies in all slots**
 - the **dispersion of the tokens**
 - the **distribution/entropy of the token frequencies**
 - the **frequencies and association strengths of elements to the slot and other slots**
 - all of that **sense-specific** (cf. Roland & Jurafsky 2002; Coleman & Bernolet 2012) and in a **probabilistic network of constructions** (cf. Roland, Dick, & Elman 2007)

To conclude, this is what we agree on and what our agenda looks like ...

- All the above must not distract from the fact that I **strongly agree with most of** what **Bybee** (and of course many others) say(s): of course,
 - usage-/**exemplar-based approaches** are on the right track
 - **domain-general mechanisms** such as analogy/similarity, chunking, and frequencies of exposure/processing are key
- but the devil lies in the detail (duh!) and we must be very careful to
 - not throw too many babies out with the bathwater and prematurely choose, or abandon, particular kinds of data and methods
 - not explore the true nature/structure of our data especially since ... it doesn't get much messier / more chaotic than linguistic data :-|

Thank you!

<http://tinyurl.com/stgries>