

- Slobin, D. I. (ed.) (1985), *The Cross-linguistic Study of Language Acquisition*, vol. 1–2. Hillsdale: Lawrence Erlbaum.
- Slobin, D. I. (ed.) (1992), *The Cross-linguistic Study of Language Acquisition*, vol. 3. Hillsdale: Lawrence Erlbaum.
- Slobin, D. I. (ed.) (1997), *The Cross-linguistic Study of Language Acquisition*, vol. 4–5. Hillsdale: Lawrence Erlbaum.
- Stern, C./Stern, W. (1907), *Die Kindersprache: Eine psychologische und sprachtheoretische Untersuchung*. Leipzig: Barth.
- Tavakolian, S. L. 1977. Structural Principles in the Acquisition of Complex Sentences. Ph.D. dissertation, University of Massachusetts.
- Thompson, S. (2002), ‘Object Complements’ and Conversation: Towards a Realistic Account. In: *Studies in Language* 26, 125–164.
- Tomasello, M./Stahl, D. (2004), Sampling Children’s Spontaneous Speech: How Much is Enough? In: *Journal of Child Language* 31, 101–121.
- Verhagen, A. (2005), *Constructions of Intersubjectivity: Discourse, Syntax, and Cognition*. Oxford: Oxford University Press.
- Vihman, M. M. (1996), *Phonological Development: The Origins of Language in the Child*. Oxford: Blackwell.
- Yip, V./Matthews, S. (2000), Syntactic Transfer in a Bilingual Child. In: *Bilingualism: Language and Cognition* 3, 193–208.
- Zipf, G. (1935), *The Psycho-biology of Language*. Boston: Houghton Mifflin.

Holger Diessel, Jena (Germany)

58. Corpora and collocations

1. Introduction
2. What are collocations?
3. Cooccurrence and frequency counts
4. Simple association measures
5. Statistical association measures
6. Finding the right measure
7. Summary and conclusion
8. Literature

1. Introduction

1.1. The controversy around collocations

The concept of *collocations* is certainly one of the most controversial notions in linguistics, even though it is based on a compelling, widely shared intuition that certain words have a tendency to occur near each other in natural language. Examples of such collocations are *cow* and *milk*, *day* and *night*, *ring* and *bell*, or the infamous *kick* and *bucket*. Other words, like *know* and *glass* or *door* and *year*, do not seem to be particularly attracted to each other. J. R. Firth (1957) introduced the term “collocations” for charac-

teristic and frequently recurrent word combinations, arguing that the meaning and usage of a word (the *node*) can to some extent be characterised by its most typical *collocates*: “You shall know a word by the company it keeps” (Firth 1957, 179). Firth was clearly aware of the limitations of this approach. He understood collocations as a convenient first approximation to meaning at a purely lexical level that can easily be operationalised (cf. Firth 1957, 181). Collocations in this Firthian sense can also be interpreted as empirical statements about the predictability of word combinations: they quantify the “mutual expectancy” (Firth 1957, 181) between words and the statistical influence a word exerts on its neighbourhood. Firth’s definition of the term remained vague, though, and it was only formalised and implemented after his death, by a group of British linguists often referred to as the Neo-Firthian school. Collocations have found widespread application in computational lexicography (Sinclair 1966, 1991), resulting in corpus-based dictionaries such as COBUILD (Sinclair 1995; see also article 8).

In parallel to the development of the Neo-Firthian school, the term “collocations” came to be used in the field of phraseology for semi-compositional and lexically determined word combinations such as *stiff drink* (with a special meaning of *stiff* restricted to a particular set of nouns), *heavy smoker* (where *heavy* is the only acceptable intensifier for *smoker*), *give a talk* (rather than *make or hold*) and *a school of fish* (rather than *group, swarm or flock*). This view has been advanced forcefully by Hausmann (1989) and has found increasingly widespread acceptance in recent years (e.g. Grossmann/Tutin 2003). It is notoriously difficult to give a rigorous definition of collocations in the phraseological sense and differentiate them from restricted word senses (most dictionaries have separate subentries for the special meanings of *stiff*, *heavy* and *school* in the examples above). There is considerable overlap between the phraseological notion of collocations and the more general empirical notion put forward by Firth (cf. the examples given above), but they are also different in many respects (e.g., *good* and *time* are strongly collocated in the empirical sense, but *a good time* can hardly be understood as a non-compositional or lexically restricted expression). This poor alignment between two interpretations of the same term has resulted in frequent misunderstandings and has led to enormous confusion on both sides. The situation is further complicated by a third meaning of “collocations” in the field of computational linguistics, where it is often used as a generic term for any lexicalised word combination that has idiosyncratic semantic or syntactic properties and may therefore require special treatment in a machine-readable dictionary or natural language processing system. This usage seems to originate with Choueka (1988) and can be found in standard textbooks, where collocations are often defined in terms of non-compositionality, non-modifiability and non-substitutability (Manning/Schütze 1999, 184). It has recently been superseded by the less ambiguous term *multiword expression* (cf. Sag et al. 2002).

An excellent overview of the competing definitions of collocations and their historical development is given by Bartsch (2004). Interestingly, she takes a middle road with her working definition of collocations as “lexically and/or pragmatically constrained recurrent co-occurrences of at least two lexical items which are in a direct syntactic relation with each other” (Bartsch 2004, 76). For a compact summary, refer to Williams (2003).

1.2. Definitions and recommended terminology

In order to avoid further confusion, a consistent terminology should be adopted. Its most important goal is to draw a clear distinction between (i) the *empirical* concept of

recurrent and predictable word combinations, which are a directly observable property of natural language, and (ii) the *theoretical* concept of lexicalised, idiosyncratic multiword expressions, defined by linguistic tests and speaker intuitions. In this article, the term “*collocations*” is used exclusively in its empirical Firthian sense (i), and we may occasionally speak of “*empirical collocations*” to draw attention to this fact. Lexicalised word combinations as a theoretical, phraseological notion (ii) are denoted by the generic term “*multiword expressions*”, following its newly established usage in the field of computational linguistics. In phraseological theory, multiword expressions are divided into sub-categories ranging from completely opaque idioms to semantically compositional word combinations, which are merely subject to arbitrary lexical restrictions (*brush teeth* rather than *scrub teeth*) or carry strong pragmatic connotations (*red rose*). A particularly interesting category in the middle of this cline are semi-compositional expressions, in which one of the words is lexically determined and has a modified or bleached meaning (classic examples are *heavy smoker* and *give a talk*). They correspond to the narrow phraseological meaning of the term “*collocations*” (cf. Grossmann/Tutin 2003) and can be referred to as “*lexical collocations*”, following Krenn (2000). As has been pointed out above, it is difficult to give a precise definition of lexical collocations and to differentiate them e. g. from specialised word senses. Because of this fuzziness and the fact that many empirical collocations are neither completely opaque nor fully compositional, similar to lexical collocations, the two concepts are easily and frequently confused.

This article is concerned exclusively with empirical collocations, since they constitute one of the fundamental notions of corpus linguistics and, unlike lexicalisation phenomena, can directly be observed in corpora. It is beyond the scope of this text to delve into the voluminous theoretical literature on multiword expressions, but see e. g. Bartsch (2004) and Grossmann/Tutin (2003) for useful pointers. There is a close connection between empirical collocations and multiword expressions, though. A thorough analysis of the collocations found in a corpus study will invariably bring up non-compositionality and lexicalisation phenomena as an explanation for many of the observed collocations (cf. the case study in section 2.2.). Conversely, theoretical research in phraseology can build on authentic examples of multiword expressions obtained from corpora, avoiding the bias of relying on introspection or stock examples like *kick the bucket* (which is a rather uncommon phrase indeed: only three instances of the idiom can be found in the 100 million words of the British National Corpus). *Multiword extraction* techniques exploit the often confusing overlap between the empirical and theoretical notions of collocation. Empirical collocations are identified as candidate multiword expressions, and then the “false positives” are weeded out by manual inspection. A more detailed account of such multiword extraction procedures can be found in section 6.2.

Following the Firthian tradition (e. g. Sinclair 1991), we define a collocation as a combination of two words that exhibit a tendency to occur near each other in natural language, i. e. to *cooccur* (but see the remarks on combinations of three or more words in section 7.1.). The term “*word pair*” is used to refer to such a combination of two words (or, more precisely, word *types*; see article 36 for the distinction between types and tokens) in a neutral way without making a commitment regarding its collocational status. In order to emphasise this view of collocations as word pairs, we will use the notation (*kick*, *bucket*) instead of e. g. *kick (the) bucket*. In general, a word pair is denoted by (w_1, w_2) , with $w_1 = \textit{kick}$ and $w_2 = \textit{bucket}$ in the previous example; w_1 and w_2 are also referred to as the *components* of the word pair. The term “word” is meant in

the widest possible sense here and may refer to surface forms, case-folded surface forms, base forms, etc. (see article 25). While collocations are most commonly understood as combinations of orthographic words, delimited by whitespace and punctuation, the concept and methodological apparatus can equally well be applied to combinations of linguistic units at other levels, ranging from morphemes to phrases and syntactic constructions (cf. article 43).

In order to operationalise our definition of collocations, we need to specify the precise circumstances under which two words can be said to “cooccur”. We also need a formal definition of the “attraction” between words reflected by their repeated cooccurrence, and a quantitative measure for the strength of this attraction. The cooccurrence of words can be defined in many different ways. The most common approaches are (i) *surface cooccurrence*, where words are said to cooccur if they appear close to each other in running text, measured by the number of intervening word tokens; (ii) *textual cooccurrence* of words in the same sentence, clause, paragraph, document, etc.; and (iii) *syntactic cooccurrence* between words in a (direct or indirect) syntactic relation, such as a noun and its modifying adjective (which tend to be adjacent in most European languages) or a verb and its object noun (which may be far apart at the surface, cf. Goldman/Nerima/Wehrli (2001, 62) for French). These three definitions of cooccurrence are described in more detail in section 3, together with appropriate methods for the calculation of cooccurrence frequency data.

The hallmark of an attraction between words is their frequent cooccurrence, and collocations are sometimes defined simply as “recurrent cooccurrences” (Smadja 1993, 147; Bartsch 2004, 11). Strictly speaking, any pair of words that cooccur at least twice in a corpus is a potential collocation according to this view. It is common to apply higher *frequency thresholds*, however, such as a minimum of 3, 5 or even 10 cooccurrences. Evert (2004, chapter 4) gives a mathematical justification for this approach (see also section 7.1.), but a more practical reason is to reduce the enormous amounts of data that have to be processed. It is not uncommon to find more than a million recurrent word pairs ($f \geq 2$) in a corpus containing several hundred million running words, but only a small proportion of them will pass a frequency threshold of $f \geq 10$ or higher, as a consequence of Zipf’s law (cf. article 37). In the following, we use the term “*recurrent word pair*” for a potential collocation that has passed the chosen frequency threshold in a given corpus.

Mere recurrence is not a sufficient indicator for a strong attraction between words, though, as will be illustrated in section 4.1. An additional measure of attraction strength is therefore needed in order to identify “true collocations” among the recurrent word pairs, or to distinguish between “strong” and “weak” collocations. The desire to generalise from recurrent word pairs in a particular corpus (as a sample of language) to collocations in the full language or sublanguage, excluding word pairs whose recurrence may be an accident of the sampling process, has led researchers to the concept of *statistical association* (Sinclair 1966, 418). Note that this mathematical meaning of “association” describes a statistical attraction between certain events and must not be confused with psychological association (as e. g. in word association norms, which have no direct connection to the statistical association between words that is of interest here). By interpreting occurrences of words as events, statistical *association measures* can be used to quantify the attraction between cooccurring words. They thus complete the formal definition of empirical collocations.

The most important association measures will be introduced in sections 4 and 5, but many other measures have been suggested in the mathematical literature and in collocation studies. Such measures assign an *association score* to each word pair, with high scores indicating strong attraction and low scores indicating weak attraction (or even repulsion) between the component words. Association scores can then be used to select “true collocations” by setting a threshold value, or to rank the set of recurrent word pairs according to the strength of their attraction (so that “strong” collocations are found at the top of the list). These uses of association scores are further explained in section 2.1. It is important to keep in mind that different association measures may lead to entirely different rankings of the word pairs (or to different sets of “true collocations”). Section 6 gives some guidance on how to choose a suitable measure.

1.3. Overview of the article

Section 2 describes the different uses of association scores and illustrates the linguistic properties of empirical collocations with a case study of the English noun *bucket*. The three types of cooccurrence (surface, textual and syntactic) are defined and compared in section 3, and the calculation of cooccurrence frequency data is explained with the help of toy examples. Section 4 introduces the concepts of statistical association and independence underlying all association measures. It also presents a selection of simple measures, which are based on a comparison of observed and expected cooccurrence frequency. Section 5 introduces more complex statistical measures based on full-fledged contingency tables. The difficulty of choosing between the large number of available measures is the topic of section 6, which discusses various methods for the comparison of association measures. Finally, section 7 addresses some open questions and extensions that are beyond the scope of this article, and lists references for further reading.

Readers in a hurry may want to start with the “executive summaries” in section 4.3. and at the beginning of section 7, which give a compact overview of the collocation identification process with simple association measures. You should also skim the examples in section 3 to understand how appropriate cooccurrence frequency data are obtained from a corpus, find out in section 4.1. how to calculate the observed cooccurrence frequency O and expected frequency E , and refer to Figure 58.4 for the precise equations of various simple association measures.

2. What are collocations?

2.1. Using association scores

Association scores as a quantitative measure of the attraction between words play a crucial role in the operationalisation of empirical collocations, next to the formal definition of cooccurrence and the appropriate calculation of cooccurrence frequency data. While the interpretation of association scores seems straightforward (high scores indicate strong attraction), they can be used in different ways to identify collocations among the recurrent word pairs found in a corpus. The first contrast to be made is whether colloca-

tivity is treated as a categorical phenomenon or as a cline, leading either to *threshold* approaches (which attempt to identify “true collocations”) or to *ranking* approaches (which place word pairs on a scale of collocational strength without strict separation into collocations and non-collocations). A second contrast concerns the grouping of collocations: the *unit* view is interested in the most strongly collocated word pairs, which are seen as independent units; the *node-collocate* view focuses on the collocates of a given node word, i.e. “the company it keeps”. The two contrasts are independent of each other in principle, although the *node-collocate* view is typically combined with a ranking approach.

In a threshold approach, recurrent word pairs whose association score exceeds a (more or less arbitrary) threshold value specified by the researcher are accepted as “true collocations”. We will sometimes refer to them as an *acceptance set* for a given association measure and threshold value. In the alternative approach, all word pairs are ranked according to their association scores. Pairs at the top of the ranked list are then considered “more collocational”, while the ones at the bottom are seen as “less collocational”. However, no categorical distinction between collocations and non-collocations is made in this approach. A third strategy combines the ranking and threshold approaches by accepting the first n word pairs from the ranked list as collocations, with n either determined interactively by the researcher or dictated by the practical requirements of an application. Typical choices are $n = 100$, $n = 500$, $n = 1000$ and $n = 2000$. Such *n-best lists* can be interpreted as acceptance sets for a threshold value determined from the corpus data (such that exactly n word pairs are accepted) rather than chosen at will. Because of the arbitrariness of pre-specified threshold values and the lack of good theoretical motivations (cf. section 4.2.), n -best lists should always be preferred over threshold-based acceptance sets. It is worth pointing out that in either case the ranking, n -best list or acceptance set depends critically on the particular association measure that has been used. The n -best lists shown in Tables 58.2 and 58.3 in section 4.3. are striking examples of this fact.

The unit view interprets collocations as pairs of words that show a strong mutual attraction, or “mutual expectancy” (Firth 1957, 181). It is particularly suitable and popular for multiword extraction tasks, where n -best lists containing the most strongly associated word pairs in a corpus are taken as candidate multiword expressions. Such candidate lists serve e.g. as base material for dictionary updates, as terminological resources for translators and technical writers, and for the semi-automatic compilation of lexical resources for natural language processing systems (e.g. Heid et al. 2000). The node-collocate view, on the other hand, focuses on the predictability of word combinations, i.e. on how a word (the node) determines its “company” (the collocates). It is well suited for the linguistic description of word meaning and usage in the Firthian tradition, where a node word is characterised by ranked lists of its collocates (Firth 1957). Following Firth (1957, 195–196) and Sinclair (1966), this view has also found wide acceptance in modern corpus-based lexicography (e.g. Sinclair 1991; Kilgariff et al. 2004), in particular for learner dictionaries such as *COBUILD* (Sinclair 1995) and the *Oxford Collocations Dictionary* (Lea 2002).

In addition to their “classic” applications in language description, corpus-based lexicography and multiword extraction, collocations and association scores have many practical uses in computational linguistics and related fields. Well-known examples include the construction of machine-readable dictionaries for machine translation and natural

language generation systems, the improvement of statistical language models, and the use of association scores as features in vector space models of distributional semantics. See Evert (2004, 23–27) for an overview and comprehensive references.

2.2. Collocations as a linguistic epiphenomenon

The goal of this section is to help readers reach an intuitive understanding of the empirical phenomenon of collocations and their linguistic properties. First and foremost, collocations are observable facts about language, i.e. primary data. From a strictly data-driven perspective, they can be interpreted as empirical predictions about the neighbourhood of a word. For instance, a verb accompanying the noun *kiss* is likely to be either *give*, *drop*, *plant*, *press*, *steal*, *return*, *deepen*, *blow* or *want*. From the explanatory perspective of theoretical linguistics, on the other hand, collocations are best characterised as an *epiphenomenon*: idioms, lexical collocations, clichés, cultural stereotypes, semantic compatibility and many other factors are hidden causes that result in the observed associations between words. In order to gain a better understanding of collocations both as an empirical phenomenon and as an epiphenomenon, we will now take a look at a concrete example, viz. how the noun *bucket* is characterised by its collocates in the British National Corpus (BNC, Aston/Burnard 1998). The data presented here are based on surface cooccurrence with a span size of 5 words, delimited by sentence boundaries (see section 3). Observed and expected frequencies were calculated as described in section 4.1. Collocates were lemmatised, and punctuation, symbols and numbers were excluded. Association scores were calculated for the measures MI and simple-II (see section 4.2.).

A first observation is that different association measures will produce entirely different rankings of the collocates. For the MI measure, the top collocates are *fourteen-record*, *ten-record*, *full-track*, *single-record*, *randomize*, *galvanized*, *groundbait*, *slop*, *spade*, *Nessie*. Most of them are infrequent words with low cooccurrence frequency (e.g., *groundbait* occurs only 29 times in the BNC). Interestingly, the first five collocates belong to a technical sense of *bucket* as a data structure in computer science; others such as *groundbait* and *Nessie* (the name of a character in the novel *Worlds Apart*, BNC file ATE) are purely accidental combinations. By contrast, the top collocates according to the simple-II measure are dominated by high-frequency cooccurrences with very common words, including several function words: *water*, *a*, *spade*, *plastic*, *size*, *slop*, *mop*, *throw*, *fill*, *with*.

A clearer picture emerges when different parts of speech among the collocates (e.g. nouns, verbs and adjectives) are listed separately, as shown in Table 58.1 for the simple-II measure. Ideally, a further distinction should be made according to the syntactic relation between node and collocate (node as subject/object of verb, prenominal adjective modifying the node, head of postnominal *of*-NP, etc.), similar to the lexicographic *word sketches* of Kilgarriff et al. (2004). Parts of speech provide a convenient approximation that does not require sophisticated automatic language processing tools. A closer inspection of the lists in Table 58.1 underlines the status of collocations as an epiphenomenon, revealing many different causes that contribute to the observed associations:

Tab. 58.1: Collocates of *bucket* in the BNC (nouns, verbs and adjectives)

noun	f	simple-ll	verb	f	simple-ll	adjective	f	simple-ll
<i>water</i>	183	1063.90	<i>throw</i>	36	165.32	<i>large</i>	37	92.72
<i>spade</i>	31	338.21	<i>fill</i>	29	129.69	<i>single-record</i>	5	79.56
<i>plastic</i>	36	242.63	<i>randomize</i>	9	115.33	<i>cold</i>	13	52.63
<i>slop</i>	14	197.65	<i>empty</i>	14	106.51	<i>galvanized</i>	4	52.35
<i>size</i>	41	193.22	<i>tip</i>	10	62.65	<i>ten-record</i>	3	49.75
<i>mop</i>	16	183.97	<i>kick</i>	12	59.12	<i>full</i>	20	46.34
<i>record</i>	38	155.64	<i>hold</i>	31	58.52	<i>empty</i>	9	36.41
<i>bucket</i>	18	138.70	<i>carry</i>	26	55.68	<i>steaming</i>	4	36.37
<i>ice</i>	22	131.68	<i>put</i>	36	48.69	<i>full-track</i>	2	33.17
<i>seat</i>	20	78.35	<i>chuck</i>	7	48.40	<i>multi-record</i>	2	33.17
<i>coal</i>	16	76.44	<i>weep</i>	7	44.14	<i>small</i>	21	30.90
<i>density</i>	11	66.78	<i>pour</i>	9	39.35	<i>leaky</i>	3	30.14
<i>brigade</i>	10	66.78	<i>douse</i>	4	37.85	<i>bottomless</i>	3	29.04
<i>algorithm</i>	9	66.54	<i>fetch</i>	7	35.22	<i>galvanised</i>	3	28.34
<i>shovel</i>	7	64.53	<i>store</i>	7	30.77	<i>iced</i>	3	25.46
<i>container</i>	10	62.40	<i>drop</i>	9	21.76	<i>clean</i>	7	25.17
<i>oats</i>	7	62.32	<i>pick</i>	11	21.74	<i>wooden</i>	6	24.14
<i>sand</i>	12	61.91	<i>use</i>	31	20.93	<i>old</i>	19	18.83
<i>Rhino</i>	7	60.50	<i>tire</i>	3	20.58	<i>ice-cold</i>	2	17.66
<i>champagne</i>	10	59.28	<i>rinse</i>	3	20.19	<i>anti-sweat</i>	1	16.58

- the well-known idiom *kick the bucket*, although many of the cooccurrences represent a literal reading of the phrase (e. g. *It was as if God had kicked a bucket of water over.*, G0P: 2750);
- proper names such as *Rhino Bucket*, a hard rock band founded in 1987;
- both lexicalised and productively formed compound nouns: *slop bucket*, *bucket seat*, *coal bucket*, *champagne bucket* and *bucket shop* (the 23rd noun collocate);
- lexical collocations like *weep buckets*, where *buckets* has lost its regular meaning and acts as an intensifier for the verb;
- cultural stereotypes and institutionalised phrases such as *bucket and spade* (which people prototypically take along when they go to a beach, even though the phrase has fully compositional meaning);
- reflections of semantic compatibility: *throw*, *carry*, *kick*, *tip*, *take*, *fetch* are typical things one can do with a bucket, and *full*, *empty*, *leaky* are some of its typical properties (or states);
- semantically similar terms (*shovel*, *mop*) and hypernyms (*container*);
- facts of life, which do not have special linguistic properties but are frequent simply because they describe a situation that often arises in the real world; a prototypical example is *bucket of water*, the most frequent noun collocate in Table 58.1;
- linguistic relevance: it is more important to talk about *full*, *empty* and *leaky* buckets than e. g. about a rusty or yellow bucket; interestingly, *old bucket* (*f* = 19) is much more frequent than *new bucket* (*f* = 3, not shown); and
- “indirect” collocates (e. g. *a bucket of cold*, *warm*, *hot*, *iced*, *steaming water*), describing typical properties of the liquid contained in a bucket.

Obviously, there are entirely different sets of collocates for each sense of the node word, which are overlaid in Table 58.1. As Firth put it: “there are the specific contrastive collocations for *light/dark* and *light/heavy*” (Firth 1957, 181). In the case of *bucket*, a technical meaning, referring to a specific data structure in computer science, is conspicuous and accounts for a considerable proportion of the collocations (*bucket brigade algorithm*, *bucket size*, *randomize to a bucket*, *store records in bucket*, *single-record bucket*,

ten-record bucket). In order to separate collocations for different word senses automatically, a sense-tagged corpus would be necessary (cf. article 26).

Observant readers may have noticed that the list of collocations in Table 58.1 is quite similar to the entry for *bucket* in the *Oxford Collocations Dictionary* (OCD, Lea 2002). This is not as surprising as it may seem at first, since the OCD is also based on the British National Corpus as its main source of corpus data (Lea 2002, viii). Obviously, collocations were identified with a technique similar to the one used here.

3. Cooccurrence and frequency counts

As has already been stated in section 1.2., the operationalisation of collocations requires a precise definition of the cooccurrence, or “nearness”, of two words (or, more precisely, word *tokens*). Based on this definition, cooccurrence frequency data for each recurrent word pair (or, more precisely, pair of word *types*) can be obtained from a corpus. Association scores as a measure of attraction between words are then calculated from these frequency data. It will be shown in section 4.1. that *cooccurrence frequency* alone is not sufficient to quantify the strength of attraction. It is also necessary to consider the occurrence frequencies of the individual words, known as *marginal frequencies*, in order to assess whether the observed cooccurrences might have come about by chance. In addition, a measure of corpus size is needed to interpret absolute frequency counts. This measure is referred to as *sample size*, following statistical terminology.

The following notation is used in this article: O for the “observed” cooccurrence frequency in a given corpus (sometimes also denoted by f , especially when specifying frequency thresholds such as $f \geq 5$); f_1 and f_2 for the marginal frequencies of the first and second component of a word pair, respectively; and N for the sample size. These four numbers provide the information needed to quantify the statistical association between two words, and they are called the *frequency signature* of the pair (Evert 2004, 36). Note that a separate frequency signature is computed for every recurrent word pair (w_1, w_2) in the corpus. The set of all such recurrent word pairs together with their frequency signatures is referred to as a *data set*.

Three different approaches to measuring nearness are introduced below and explained with detailed examples: *surface*, *textual* and *syntactic* cooccurrence. For each type of cooccurrence, an appropriate procedure for calculating frequency signatures (O, f_1, f_2, N) is described. The mathematical reasons behind these procedures will become clear in section 5. The aim of the present section is to clarify the logic of computing cooccurrence frequency data. Practical implementations that can be applied to large corpora use more efficient algorithms, especially for surface cooccurrences (e.g. Gil/Dias 2003; Terra/Clarke 2004).

3.1. Surface cooccurrence

The most common approach in the Firthian tradition defines cooccurrence by surface proximity, i. e. two words are said to cooccur if they appear within a certain distance or *collocational span*, measured by the number of intervening word tokens. Surface cooccur-

rence is often, though not always combined with a node-collocate view, looking for collocates within the collocational spans around the instances of a given node word.

Span size is the most important choice that has to be made by the researcher. The most common values range from 3 to 5 words (e.g. Sinclair 1991), but many other span sizes can be found in the literature. Some studies in computational linguistics have focused on bigrams of immediately adjacent words, i.e. a span size of 1 (e.g. Choueka 1988; Schone/Jurafsky 2001), while others have used span sizes of dozens or hundreds of words, especially in the context of distributional semantics (Schütze 1998). Other decisions are whether to count only word tokens or all tokens (including punctuation and numbers), how to deal with multiword units (does *out of* count as a single token or as two tokens?), and whether cooccurrences are allowed to cross sentence boundaries.

A vast deal of coolness and a peculiar degree of judgement, are requisite in catching a **hat**. A man must not be precipitate, or he runs over it; he must not rush into the opposite extreme, or he loses it altogether. [...] There was a fine gentle wind, and Mr. Pickwick's **hat** rolled sportively before it. The wind puffed, and Mr. Pickwick puffed, and the **hat** rolled over and over, as merrily as a lively porpoise in a strong tide; and on it might have *rolled*, far beyond Mr. Pickwick's reach, had not its course been providentially stopped, just as that gentleman was on the point of resigning it to its fate.

Fig. 58.1: Illustration of surface cooccurrence for the word pair (*hat*, *roll*)

Figure 58.1 shows surface cooccurrences between the words *hat* (in bold face, as node) and *roll* (in italics, as collocate). The span size is 4 words, excluding punctuation and limited by sentence boundaries. Collocational spans around instances of the node word *hat* are indicated by brackets below the text. There are two cooccurrences in this example, in the second and third span, hence $O = 2$. Note that multiple instances of a word in the same span count as multiple cooccurrences, so for *hat* and *over* we would also calculate $O = 2$ (with both cooccurrences in the third span). The marginal frequencies of the two words are given by their overall occurrence counts in the text, i.e. $f_1 = 3$ for *hat* and $f_2 = 3$ for *roll*. The sample size N is simply the total number of tokens in the corpus, counting only tokens that are relevant to the definition of spans. In our example, N is the number of word tokens excluding punctuation, i.e. $N = 111$ for the text shown in Figure 58.1. If we include punctuation tokens in our distance measurements, the sample size would accordingly be increased to $N = 126$ (9 commas, 4 full stops and 2 semicolons). The complete frequency signature for the pair (*hat*, *roll*) is thus (2,3,3,111). Of course, realistic data will have much larger sample sizes, and the marginal frequencies are usually considerably higher than the cooccurrence frequency.

Collocational spans can also be asymmetric, and are generally written in the form (L_k , R_n) for a span of k tokens to the left of the node word and n tokens to its right. The symmetric spans in the example above would be described as (L_4 , R_4). Asymmetric spans introduce an asymmetry between node word and collocate that is absent from most other approaches to collocations. For a one-sided span (L_4 , R_0) to the left of the node word, there would be 2 cooccurrences of the pair (*roll*, *hat*) in Figure 58.1, but none of the pair (*hat*, *roll*). A special case are spans of the form (L_0 , R_1), where cooccurrences are ordered pairs of immediately adjacent words, often referred to as bigrams in computational linguistics. Thus, *took place* would be a bigram cooccurrence of the lemma pair (*take*, *place*), but neither *place taken* nor *take his place* would count as cooccurrences.

3.2. Textual cooccurrence

A second approach considers words to cooccur if they appear in the same textual unit. Typically, such units are sentences or utterances, but with the recent popularity of Google searches and the Web as corpus (see article 18), cooccurrence within (Web) documents has found more widespread use.

One criticism against surface cooccurrence is the arbitrary choice of the span size. For a span size of 3, *throw a birthday party* would be accepted as a cooccurrence of (*throw*, *party*), but *throw a huge birthday party* would not. This is particularly counterintuitive for languages with relatively free word order, where closely related words can be far apart on the surface. In such languages, textual cooccurrence within the same sentence may provide a more appropriate collocational span. Textual cooccurrence also captures weaker dependencies, in particular those caused by paradigmatic semantic relations. For example, if an English sentence contains the noun *bucket*, it is quite likely to contain the noun *mop* as well (although the connection is far weaker than for *water* or *spade*), but the two nouns will not necessarily be near each other in the sentence.

A vast deal of coolness and a peculiar degree of judgement, are requisite in catching a <u>hat</u> .	hat	—
A man must not be precipitate, or he runs <i>over</i> it ;	—	over
he must not rush into the opposite extreme, or he loses it altogether.	—	—
There was a fine gentle wind, and Mr. Pickwick's <u>hat</u> rolled sportively before it.	hat	—
The wind puffed, and Mr. Pickwick puffed, and the <u>hat</u> rolled <i>over</i> and <i>over</i> as merrily as a lively porpoise in a strong tide ;	hat	over

Fig. 58.2: Illustration of textual cooccurrence for the word pair (*hat*, *over*)

The definition of textual cooccurrence and the appropriate procedure for computing frequency signatures are illustrated in Figure 58.2, for the word pair (*hat*, *over*) and sentences as textual segments. There is one cooccurrence of *hat* and *over* in the last sentence of this text sample, hence $O = 1$. In contrast to surface cooccurrence, the count is 1 even though there are two instances of *over* in the sentence. Similarly, the marginal frequencies are given by the number of sentences containing each word, ignoring multiple occurrences in the same sentence: hence $f_1 = 3$ and $f_2 = 2$ (although there are three instances each of *hat* and *over* in the text sample). The sample size $N = 5$ is the number of sentences in this case. The complete frequency signature of (*hat*, *over*) is thus (1,3,2,5), whereas for surface cooccurrence within the spans shown in Figure 58.1 it would have been (2,3,3,79).

3.3. Syntactic cooccurrence

In this more restrictive approach, words are only considered to be near each other if there is a direct syntactic relation between them. Examples are a verb and its object (or subject) noun, prenominal adjectives (in English and German) and nominal modifiers

(the pattern N of N in English, genitive noun phrases in German). Sometimes, indirect relations might also be of interest, e.g. a verb and the adjectival modifier of its object noun, or a noun and the adjective modifying a postnominal *of*-NP. The latter pattern accounts for several surface collocations of the noun *bucket* such as *a bucket of iced, cold, steaming water* (cf. Table 58.1). Collocations for different types of syntactic relations are usually treated separately. From a given corpus, one might extract a data set of verbs and their object nouns, another data set of verbs and subject nouns, a data set of adjectives modifying nouns, etc. Syntactic cooccurrence is particularly appropriate if there may be long-distance dependencies between collocates: unlike surface cooccurrence, it does not set an arbitrary distance limit, but at the same time it does not introduce as much “noise” as textual cooccurrence. Syntactic cooccurrence is often used for multiword extraction, since many types of lexicalised multiword expressions tend to appear in specific syntactic patterns such as verb + object noun, adjective + noun, adverb + verb, verb + predicated adjective, delexical verb + noun, etc. (see Bartsch 2004, 11).

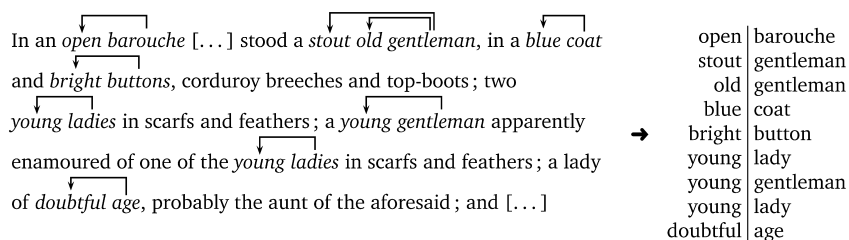


Fig. 58.3: Illustration of syntactic cooccurrence (nouns modified by prenominal adjectives)

Frequency signatures for syntactic cooccurrence are obtained in a more indirect way, illustrated in Figure 58.3. First, all instances of the desired syntactic relation are identified, in this case modification of nouns by prenominal adjectives. Then the corresponding arguments are compiled into a list with one entry for each instance of the syntactic relation (shown on the right of Figure 58.3). Note that the list entries are lemmatised here, but e.g. case-folded word forms could have been used as well. Just as the original corpus is understood as a sample of language, the list items constitute a sample of the targeted syntactic relation, and Evert (2004) refers to them as “pair tokens”. Cooccurrence frequency data are computed from this sample, while all word tokens that do not occur in the relation of interest are disregarded. For the word pair (*young*, *gentleman*), we find one cooccurrence in the list of pair tokens, i.e. $O = 1$. The marginal frequencies are given by the total numbers of entries containing one of the component words, $f_1 = 3$ and $f_2 = 3$, and the sample size is the total number of list entries, $N = 9$. The frequency signature of (*young*, *gentleman*) as a syntactic adjective-noun cooccurrence is thus (1,3,3,9).

3.4. Comparison

Collocations according to surface cooccurrence have proven useful in corpus-linguistic and lexicographic research (cf. Sinclair 1991). They strike a balance between the restricted notion of syntactic cooccurrence (esp. when only a single type of syntactic rela-

tion is considered) and the very broad notion of textual cooccurrence. The number of recurrent word pairs extracted from a corpus is also more manageable than for textual cooccurrence. In this respect, syntactic cooccurrence are even more practical. A popular application of surface cooccurrence in computational linguistics are word space models of distributional semantics (Schütze 1998; Sahlgren 2006). As an alternative to the surface approach, Kilgarriff et al. (2004) collect syntactic collocates from different types of syntactic relations and display them as a *word sketch* of the node word.

Textual cooccurrence is easier to implement than surface cooccurrence, and more robust against certain types of non-randomness such as term clustering, especially when the textual units used are entire documents (cf. the discussion of non-randomness in article 36). However, it tends to create huge data sets of recurrent word pairs that can be challenging even for powerful modern computers.

Syntactic cooccurrence separates collocations of different syntactic types, which are overlaid in frequency data according to surface cooccurrence, and discards many indirect and accidental cooccurrences. It should thus be easier to find suitable association measures to quantify the collocativity of word pairs. Evert (2004, 19) speculates that different measures might be appropriate for different types of syntactic relations. Syntactic cooccurrence is arguably most useful for the identification of multiword expressions, which are typically categorised according to their syntactic structure. However, it requires an accurate syntactic analysis of the source corpus, which will have to be performed with automatic tools in most cases. For prenominal adjectives, the analysis is fairly easy in English and German (Evert/Kermes 2003), while for German verb-object relations, it is extremely difficult to achieve satisfactory results: recent syntactic parsers achieve dependency F-scores of 70%–75% (Schiehlen 2004). Outspoken advocates of syntactic cooccurrence include Daille (1994), Goldman/Nerima/Wehrli (2001), Bartsch (2004) and Evert (2004).

Leaving such practical and philosophical considerations aside, frequency signatures computed according to the different types of cooccurrence can disagree substantially for the same word pair. For example, the frequency signatures of (*short*, *time*) in the Brown corpus are: (16,135,457,59710) for syntactic cooccurrence (of prenominal adjectives), (27,213,1600,1170811) for (L5, R5) surface cooccurrence, and (32,210,1523,52108) for textual cooccurrence within sentences.

4. Simple association measures

4.1. Expected frequency

It might seem natural to use the cooccurrence frequency O as an association measure to quantify the strength of collocativity (e.g. Choueka 1988). This is not sufficient, however; the marginal frequencies of the individual words also have to be taken into account. To illustrate this point, consider the following example. In the Brown corpus, the bigram *is to* is highly recurrent. With $O = 260$ cooccurrences it is one of the most frequent bigrams in the corpus. However, both components are frequent words themselves: *is* occurs roughly 10,000 times and *to* roughly 26,000 times among 1 million word tokens. If the words in this corpus were rearranged in completely random order, thereby remov-

ing all associations between cooccurring words, we would still expect to see the sequence *is to* approx. 260 times. The high cooccurrence frequency of *is to* therefore does not constitute evidence for a collocation; on the contrary, it indicates that *is* and *to* are not attracted to each other at all. The expected number of cooccurrences for a completely “uncollocational” word pair has been derived by the following reasoning: *to* occurs 26 times every 1,000 words on average. If there is no association between *is* and *to*, then each of the 10,000 instances of *is* in the Brown corpus has a chance of 26/1,000 to be followed by *to*. Therefore, we expect around $10,000 \times (26/1,000) = 260$ occurrences of the bigram *is to*, provided that there is indeed no association between the words. Of course, even in a perfectly randomised corpus there need not be exactly 260 cooccurrences: statistical calculations compute averages across large numbers of samples (formally called *expectations*), while the precise value in a corpus is subject to unpredictable random variation (see article 36).

The complete absence of association, as between words in a randomly shuffled corpus, is called *independence* in mathematical statistics. What we have calculated above is the *expected value* for the number of cooccurrences in a corpus of 1 million words, under the *null hypothesis* that *is* and *to* are independent. In analogy to the *observed frequency* O of a word pair, the expected value under the null hypothesis of independence is denoted E and referred to as the *expected frequency* of the word pair. Expected frequency serves as a reference point for the interpretation of O : the pair is only considered collocational if the observed cooccurrence frequency is substantially greater than the expected frequency, $O \gg E$. Using the formal notation of section 3, the marginal frequencies of (*is*, *to*) are $f_1 = 10,000$ and $f_2 = 26,000$. The sample size is $N = 1,000,000$ tokens, and the observed frequency is $O = 260$. Expected frequency is thus given by the equation $E = f_1 \cdot (f_2 / N) = f_1 f_2 / N = 260$. While the precise calculation of expected frequency is different for each type of cooccurrence, it always follows the basic scheme $f_1 f_2 / N$. For textual and syntactic cooccurrence, the standard formula $E = f_1 f_2 / N$ can be used directly. For surface cooccurrence, an additional factor k represents the total span size, i.e. $E = k f_1 f_2 / N$. This factor is $k = 10$ for a symmetric span of 5 words (L5, R5), $k = 4$ for a span (L3, R1), and $k = 1$ for simple bigrams (L0, R1).

4.2. Essential association measures

A *simple association measure* interprets observed cooccurrence frequency O by comparison with the expected frequency E , and calculates an *association score* as a quantitative measure for the attraction between two words. The most important and widely used simple association measures are shown in Figure 58.4. In the following paragraphs, their mathematical background and some important properties will be explained.

$$\begin{aligned} \text{MI} &= \log_2 \frac{O}{E} & \text{MI}^k &= \log_2 \frac{O^k}{E} & \text{local-MI} &= O \cdot \log_2 \frac{O}{E} \\ \text{z-score} &= \frac{O - E}{\sqrt{E}} & \text{t-score} &= \frac{O - E}{\sqrt{O}} & \text{simple-ll} &= 2 \left(O \cdot \log \frac{O}{E} - (O - E) \right) \end{aligned}$$

Fig. 58.4: A selection of simple association measures

The most straightforward and intuitive way to relate O and E is to use the ratio O/E as an association measure. For instance, $O/E = 10$ means that the word pair cooccurs 10 times more often than would be expected by chance, indicating a certain degree of collocativity. Since the value of O/E can become extremely high for large sample size (because $E \ll 1$ for many word pairs), it is convenient and sensible to measure association on a (base-2) logarithmic scale. This measure can also be derived from information theory, where it is interpreted as the number of bits of “shared information” between two words and known as (*pointwise*) *mutual information* or simply MI (Church/Hanks 1990, 23). A MI value of 0 bits corresponds to a word pair that cooccurs just as often as expected by chance ($O = E$); 1 bit means twice as often ($O = 2E$), 2 bits mean 4 times as often, 10 bits about 1000 times as often, etc. A negative MI value indicates that a word pair cooccurs less often than expected by chance: half as often for -1 bit, a quarter as often for -2 bits, etc. Thus, negative MI values constitute evidence for a “repulsion” between two words, the pair forming an *anti-collocation*.

The MI measure exemplifies two general *conventions for association scores* that all association measures should adhere to. (i) Higher scores indicate stronger attraction between words, i.e. a greater degree of collocativity. In particular, repulsion, i.e. $O < E$, should result in very low association scores. (ii) Ideally, an association measure should distinguish between *positive* association ($O > E$) and *negative* association ($O < E$), assigning positive and negative scores, respectively. A strong negative association would thus be indicated by a large negative value. As a consequence, the null hypothesis of independence corresponds to a score of 0 for such association measures. It is easy to see that MI satisfies both conventions: the more O exceeds E , the larger the association score will be; for $O = E$, the MI value is $\log_2 1 = 0$. Most, though not all association measures follow at least the first convention (we will shortly look at an important exception in the form of the *simple-II* measure).

In practical applications, MI was found to have a tendency to assign inflated scores to low-frequency word pairs with $E \ll 1$, especially for data from large corpora. Thus, even a single cooccurrence of two word types might result in a fairly high association score. In order to counterbalance this low-frequency bias of MI, various heuristic modifications have been suggested. The most popular one multiplies the denominator with O in order to increase the influence of observed cooccurrence frequency compared to the expected frequency, resulting in the formula $\log_2(O^2/E)$. Multiplication with O can be repeated to strengthen the counterbalancing effect, leading to an entire family of measures MI^k with $k \geq 1$, as shown in Figure 58.4. Common choices for the exponent are $k = 2$ and $k = 3$. Daille (1994) has systematically tested values $k = 2, \dots, 10$ and found $k = 3$ to work best for her purposes. An alternative way to reduce the low-frequency bias of MI is to multiply the entire formula with O , resulting in the *local-MI* measure. Unlike the purely heuristic MI^k family, local-MI can be justified by an information-theoretic argument (Evert 2004, 89) and its value can be interpreted as bits of information. Although not immediately obvious from its equation, local-MI fails to satisfy the first convention for association scores in the case of strong negative association: for fixed expected frequency E , the score reaches a minimum at $O = E/\exp(1)$, and then increases for smaller O . Local-MI distinguishes between positive and negative association, though, and satisfies both conventions if only word pairs with positive association are considered. The measures MI^k satisfy the first convention, but violate the second convention for all $k > 1$. It has been pointed out above that MI assigns high association

scores whenever O exceeds E by a large amount, even if the absolute cooccurrence frequency is as low as $O = 1$ (and $E \ll 1$). In other words, MI only looks at what is known as *effect size* in statistics and does not take into account how much *evidence* the observed data provide. We will return to the distinction between effect-size measures and evidence-based measures in section 6. Here, we introduce three simple association measures from the latter group.

A *z-score* is a standardised measure for the amount of evidence provided by a sample against a simple null hypothesis such as $O = E$ (see article 36). In our case, the general rule for calculating z-scores leads to the equation shown in Figure 58.4. Z-scores were first used by Dennis (1965, 69) as an association measure, and later by Berry-Rogghe (1973, 104). They distinguish between positive and negative association: $O > E$ leads to $z > 0$ and $O < E$ to $z < 0$. Z-scores can be interpreted by comparison with a standard normal distribution, providing theoretically motivated cut-off thresholds for the identification of “true collocations”. An absolute value $|z| > 1.96$ is generally considered sufficient to reject the null hypothesis, i.e. to provide significant evidence for a (positive or negative) association; a more conservative threshold is $|z| > 3.29$. When used as an association measure, z-score tends to yield much larger values, though, and most word pairs in a typical data set are highly significant. For instance, 80% of all distinct word bigrams in the Brown corpus have $|z| > 1.96$, and almost 70% have $|z| > 3.29$. Recent studies avoid standard thresholds and use z-scores only to rank word pairs or select n-best lists.

A fundamental problem of the z-score measure is the normal approximation used in its mathematical derivation, which is valid only for sufficiently high expected frequency E . While there is no clearly defined limit value, the approximation becomes very inaccurate if $E < 1$, which is often the case for large sample sizes (e.g., 89% of all bigrams in the Brown corpus have $E < 1$). Violation of the normality assumption leads to highly inflated z-scores and a low-frequency bias similar to the MI measure. In order to avoid this low-frequency bias, various other significance measures have been suggested, based on more “robust” statistical tests. One possibility is the *t-score* measure, which replaces E in the denominator of z-score by O . This measure has been widely used in computational lexicography following its introduction into the field by Church et al. (1991, section 2.2.). See Evert (2004, 82–83) for a criticism of its derivation from the statistical *t* test, which is entirely inappropriate for corpus frequency data.

Dunning (1993) advocated the use of likelihood-ratio tests, which are also more robust against low expected frequencies than z-score. For a simple measure comparing O and E , the likelihood-ratio procedure leads to the *simple-ll* equation in Figure 58.4. It can be shown that simple-ll scores are always non-negative and violate both conventions for association scores. Because the underlying likelihood-ratio test is a *two-sided* test, the measure does not distinguish between $O \gg E$ and $O \ll E$, assigning high positive scores in both cases. This detail is rarely mentioned in publications and textbooks and may easily be overlooked. A general procedure can be applied to convert a two-sided association measure like *simple-ll* into a one-sided measure that satisfies both conventions: association scores are calculated in the normal way and then multiplied with -1 for all word pairs with $O < E$. This procedure is applicable if association scores of the two-sided measure are always non-negative and high scores are assigned to strong negative associations. For the resulting transformed measure, significance is indicated by the absolute value of an association score, while positive and negative association are distinguished by its sign.

Similar to the z-score measure, simple-ll measures significance (i.e. the amount of evidence against the null hypothesis) on a standardised scale, known as a chi-squared distribution with one degree of freedom, or χ^2_1 for short. Theoretically motivated cut-off thresholds corresponding to those for z-scores are $|ll| > 3.84$ and $|ll| > 10.83$, but the same reservations apply: many word pairs achieve scores far above these thresholds, so that they are not a meaningful criterion for the identification of “true collocations”.

Article 36 gives detailed explanations of statistical concepts such as *significance*, *effect size*, *hypothesis test*, *one-sided* vs. *two-sided* test, *z-score* and *normal distribution* that have been used in this section.

4.3. Simple association measures in a nutshell

The preceding section has introduced a basic selection of simple association measures. These measures quantify the “attraction” between two words, i.e. their statistical association, by comparing observed cooccurrence frequency O against E , the expected frequency under the null hypothesis of independence (i.e. complete absence of association). E is important as a reference point for the interpretation of O , since two frequent words might cooccur quite often purely by chance. Most association measures follow the convention that higher association scores indicate stronger (positive) association. Many measures also differentiate between positive association ($O > E$), indicated by positive scores, and negative association ($O < E$), indicated by negative scores. Two-sided measures fail to make any distinction between positive and negative association, but can be converted into one-sided measures with an explicit test for $O > E$.

The association measures listed in Figure 58.4 offer a number of different angles on collocativity that are sufficient for many purposes. Except for the heuristic MI^k family, all measures have theoretical motivations, allowing a meaningful interpretation of the computed association scores. As has been exemplified with the standard z-score thresholds, one should not put too much weight on such interpretations, though. Cooccurrence data do not always satisfy the assumptions made by statistical hypothesis tests, and heuristic measures may be just as appropriate.

Association measures can be divided into two general groups: measures of *effect size* (MI and MI^k) and measures of *significance* (z-score, t-score and simple-ll). The former ask the question “how strongly are the words attracted to each other?” (operationalised as “how much does observed cooccurrence frequency exceed expected frequency?”), while the latter ask “how much evidence is there for a positive association between the words, no matter how small effect size is?” (operationalised as “how unlikely is the null hypothesis that the words are independent?”). The two approaches to measuring association are not entirely unrelated: a word pair with large “true” effect size is also more likely to show significant evidence against the null hypothesis in a sample. However, there is an important difference between the two groups. Effect-size measures typically fail to account for sampling variation and are prone to a low-frequency bias (small E easily leads to spuriously high effect size estimates, even for $O = 1$ or $O = 2$), while significance measures are often prone to a high-frequency bias (if O is sufficiently large, even a small relative difference between O and E , i.e. a small effect size, can be highly significant).

Of the significance measures shown in Figure 58.4, *simple-ll* is the most accurate and robust choice. Z-score has a strong low-frequency bias because the approximations used in its derivation are not valid for $E < 1$, while t-score has been derived from an inappropriate hypothesis test. Nonetheless, t-score has proven useful for certain applications, especially the identification of certain types of multiword expressions (see section 6.2.). It has to be kept in mind that *simple-ll* is a two-sided measure and assigns high scores both to positive and negative associations. If only positive associations are of interest (as is the case for most studies), then word pairs with $O < E$ should be discarded. Alternatively, *simple-ll* can be transformed into a one-sided measure that satisfies both conventions for association scores (by multiplying scores with -1 if a word pair has $O < E$).

Association measures with a background in information theory take a different approach, which at first sight seems appropriate for the interpretation of collocations as mutually predictable word combinations (e. g. Sinclair 1966, 414). They ask the question “to what extent do the occurrences of a word w_1 determine the occurrences of another word w_2 ?”, and vice versa, based on the information-theoretic notion of mutual information (MI). Interestingly, different variants of MI lead to measures with entirely different properties: pointwise MI is a measure of effect size, while local-MI is very similar to *simple-ll* and hence has to be considered a measure of significance.

It is probably impossible to choose a single most appropriate association measure (cf. the discussion in section 6). The recommended strategy is therefore to apply *simple-ll*, t-score and MI as proven association measures with well-understood mathematical properties, in order to obtain three entirely different perspectives on the cooccurrence data. MI should always be combined with a frequency threshold to counteract its low-fre-

Tab. 58.2: Collocates of *bucket in* the BNC according to the association measures *simple-ll*, t-score, MI, and MI with frequency threshold $f \geq 5$

collocate	f	f_2	<i>simple-ll</i>
<i>water</i>	184	37012	1083.18
<i>a</i>	590	2164246	449.30
<i>spade</i>	31	465	342.31
<i>plastic</i>	36	4375	247.65
<i>size</i>	42	14448	203.36
<i>slop</i>	17	166	202.30
<i>mop</i>	20	536	197.68
<i>throw</i>	38	11308	194.66
<i>fill</i>	37	10722	191.44
<i>with</i>	196	658584	171.78

collocate	f	f_2	t-score
<i>a</i>	590	2164246	15.53
<i>water</i>	184	37012	13.30
<i>and</i>	479	2616723	10.14
<i>with</i>	196	658584	9.38
<i>of</i>	497	3040670	8.89
<i>the</i>	832	6041238	8.26
<i>into</i>	87	157565	7.67
<i>size</i>	42	14448	6.26
<i>in</i>	298	1937966	6.23
<i>record</i>	43	29404	6.12

collocate	f	f_2	MI
<i>fourteen-record</i>	4	4	13.31
<i>ten-record</i>	3	3	13.31
<i>multi-record</i>	2	2	13.31
<i>two-record</i>	2	2	13.31
<i>a-row</i>	1	1	13.31
<i>anti-sweat</i>	1	1	13.31
<i>axe-blade</i>	1	1	13.31
<i>bastarding</i>	1	1	13.31
<i>dippermouth</i>	1	1	13.31
<i>Dok</i>	1	1	13.31

collocate	$f \geq 5$	f_2	MI
<i>single-record</i>	5	8	12.63
<i>randomize</i>	10	57	10.80
<i>slop</i>	17	166	10.03
<i>spade</i>	31	465	9.41
<i>mop</i>	20	536	8.57
<i>oats</i>	7	286	7.96
<i>shovel</i>	8	358	7.83
<i>rhino</i>	7	326	7.77
<i>synonym</i>	7	363	7.62
<i>bucket</i>	18	1356	7.08

Tab. 58.3: Most strongly collocated **bigrams in** the Brown corpus according to the association measures simple-ll, t-score, MI with frequency threshold $f \geq 10$, and MI with frequency threshold $f \geq 50$

bigram	$f \geq 10$	f_1	f_2	simple-ll
<i>of the</i>	9702	34036	58451	13879.8
<i>in the</i>	6018	19615	58451	9302.3
<i>it is</i>	1482	8409	9415	5612.9
<i>on the</i>	2459	5990	58451	4972.9
<i>United States</i>	395	480	600	4842.6
<i>it was</i>	1338	8409	9339	4831.2
<i>to be</i>	1715	25106	6275	4781.1
<i>had been</i>	760	5107	2460	4599.8
<i>have been</i>	650	3884	2460	4084.0
<i>has been</i>	567	2407	2460	3944.9

bigram	$f \geq 10$	f_1	f_2	t-score
<i>of the</i>	9702	34036	58451	76.30
<i>in the</i>	6018	19615	58451	61.33
<i>on the</i>	2459	5990	58451	41.83
<i>to be</i>	1715	25106	6275	37.23
<i>it is</i>	1482	8409	9415	36.24
<i>it was</i>	1338	8409	9339	34.22
<i>at the</i>	1654	5032	58451	32.72
<i>to the</i>	3478	25106	58451	31.62
<i>from the</i>	1410	4024	58451	30.66
<i>he was</i>	1110	9740	9339	30.32

bigram	$f \geq 10$	f_1	f_2	MI
<i>Hong Kong</i>	11	11	11	16.34
<i>gon na</i>	16	16	16	15.80
<i>Viet Nam</i>	14	16	14	15.80
<i>Simms Purdew</i>	12	16	12	15.80
<i>Pathet Lao</i>	10	10	17	15.71
<i>El Paso</i>	10	19	11	15.41
<i>Lo Shu</i>	21	21	21	15.40
<i>Puerto Rico</i>	21	24	21	15.21
<i>unwed mothers</i>	10	12	26	14.83
<i>carbon tetrachloride</i>	18	30	19	14.81

bigram	$f \geq 50$	f_1	f_2	MI
<i>Los Angeles</i>	50	51	50	14.12
<i>Rhode Island</i>	100	105	175	12.27
<i>Peace Corps</i>	55	171	109	11.39
<i>per cent</i>	146	371	155	11.17
<i>United States</i>	395	480	600	10.29
<i>President Kennedy</i>	54	374	156	9.72
<i>years ago</i>	138	793	246	9.33
<i>fiscal year</i>	58	118	701	9.32
<i>New York</i>	303	1598	309	9.12
<i>United Nations</i>	51	480	175	9.11

quency bias. As an example, and to illustrate the different properties of these association measures, **Table 58.2** shows the collocates of *bucket* in the British National Corpus (following the case study in section 2.2.), according to **simple-ll**, t-score, MI without frequency threshold, and MI with an additional frequency threshold of $f \geq 5$. Table 58.3 gives a second example for word bigrams in the Brown corpus (excluding punctuation). Obviously, **simple-ll and especially t-score focus on frequent grammatical patterns like *of the* or *to be***. More interesting bigrams can only be found if separate lists are generated for each part-of-speech combination. The top collocations according to MI, on the other hand, tend to be proper names and other very fixed combinations. Their cooccurrence frequency is often close to the applied frequency threshold.

5. Statistical association measures

The simple association measures introduced in section 4 are convenient and offer a range of different perspectives on collocativity. However, two serious shortcomings make this approach unsatisfactory from a theoretical point of view and may be problematic for certain types of applications. The first of these problems is most easily explained with a worked example. In a corpus of about a million words, you might find that the bigrams $A = \text{the Iliad}$ and $B = \text{must also}$ both occur $O = 10$ times, with the same expected frequency $E = 1$. Therefore, any simple measure will assign the same association score to both bigrams. However, bigram A is a combination of a very frequent word (*the* with, say, $f_1 = 100,000$) and an infrequent word (*Iliad* with $f_2 = 10$), while B combines two words of intermediate frequency (*must* and *also* with $f_1 = f_2 = 1,000$). Using the formula $E = f_1 f_2 / N$ from section 4.1., you can easily check that the expected frequency is indeed $E = 1$ for both bigrams. While O exceeds E by the same amount for *the Iliad* as for *must also*, it is intuitively obvious that bigram A is much more strongly connected than bigram

B. In particular, $O = 10$ is the highest cooccurrence frequency that can possibly be observed for these two words (since $O \leq f_1, f_2$): every instance of *Iliad* in the corpus is preceded by an instance of *the*. For bigram B, on the other hand, the words *must* and *also* could have cooccurred much more often than 10 times. One might argue that A should therefore obtain a higher association score than B, at least for certain applications.

The second limitation of simple association measures is of a more theoretical nature. We made use of statistical concepts and methods to define measures with a meaningful interpretation, but did not apply the procedures with full mathematical rigour. In statistical theory, measures of association and tests for the independence of events are always based on a cross-classification of a random sample of certain items. An appropriate representation of cooccurrence frequency data in the form of *contingency tables* is described in section 5.1., with different rules for each type of cooccurrence. Then several widely used statistical association measures are introduced in section 5.2.

We will see in section 6 that simple association measures often give close approximations to the more sophisticated association measures introduced below. Therefore, they are sufficient for many applications, so that the computational and mathematical complexities of the rigorous statistical approach can be avoided.

5.1. Contingency tables

A rigorous statistical approach to measuring association is based on contingency tables representing the cross-classification of a set of items. Such tables naturally take marginal frequencies into account, unlike a simple comparison of O against E . As a first step, we have to define the set of cooccurrence items in a meaningful way, which is different for each type of cooccurrence. Then a separate contingency table is calculated for every word pair (w_1, w_2) , using the presence of w_1 and w_2 in each cooccurrence item as factors for the cross-classification.

	w_2	$\neg w_2$					
w_1	O_{11}	O_{12}	$= R_1$	w_1	$E_{11} = \frac{R_1 C_1}{N}$	$E_{12} = \frac{R_1 C_2}{N}$	
$\neg w_1$	O_{21}	O_{22}	$= R_2$	$\neg w_1$	$E_{21} = \frac{R_2 C_1}{N}$	$E_{22} = \frac{R_2 C_2}{N}$	
	$= C_1$	$= C_2$	$= N$				

Fig. 58.5: General form of the contingency table of observed frequencies with row and column marginals (left panel), and contingency table of expected frequencies under the null hypothesis of independence (right panel)

The resulting contingency table (left panel of Figure 58.5) has four *cells*, representing the items containing both w_1 and w_2 (O_{11} , equivalent to the observed cooccurrence frequency O), the items containing w_1 but not w_2 (O_{12}), the items containing w_2 but not w_1 (O_{21}), and the items containing neither of the two words (O_{22}). These *observed frequencies* add

up to the total number of items or *sample size*, since every item has to be classified into exactly one cell of the table. The row and column sums, also called *marginal frequencies* (as they are written in the margins of the table), play an important role in the statistical analysis of contingency tables. The first row sum R_1 corresponds to the number of cooccurrence items containing w_1 , and is therefore usually equal to f_1 (except for surface cooccurrence, see below), while the first column sum C_1 is equal to f_2 . This equivalence explains the name “marginal frequencies” for f_1 and f_2 .

As in the case of simple association measures, the statistical analysis of contingency tables is based on a comparison of the observed frequencies O_{ij} with expected frequencies under the null hypothesis that the factors defining rows and columns of the table are statistically independent (which is the mathematical equivalent of the intuitive notion of independence between w_1 and w_2 introduced in section 4.1.). In contrast to the simple approach, we are not only interested in the expected number of cooccurrences of w_1 and w_2 , but have to compute expected frequencies for all four cells of the contingency table, according to the equations shown in the right panel of Figure 58.5. Note that $O_{11} = O$ and $E_{11} = E$, so statistical contingency tables are a genuine extension of the previous approach. The statistical association measures introduced in section 5.2. below are formulated in terms of observed frequencies O_{ij} and expected frequencies E_{ij} , the marginals R_i and C_j , and the sample size N . This standard notation follows Evert (2004) and allows equations to be expressed in a clean and readable form.

The definition of appropriate contingency tables is most straightforward for syntactic cooccurrence. The pair tokens on the right of Figure 58.3 can naturally be interpreted as a set of cooccurrence items. If the first word is w_1 , an item is classified into the first row of the contingency table for the pair (w_1, w_2) , otherwise it is classified into the second row. Likewise, the item is classified into the first column if the second word is w_2 and into the second column if it is not. This procedure is illustrated in the left panel of Figure 58.6. The first row sum R_1 equals the total number of cooccurrence items containing w_1 as first element, and the first column sum equals the number of items containing w_2 as second element. This corresponds to the definition of f_1 and f_2 for syntactic cooccurrence in section 3.3. The example in the right panel of Figure 58.6 shows a contingency table for the word pair (*young, gentleman*) obtained from the sample of adjective-noun cooccurrences in Figure 58.3. Since there are nine instances of adjectival modification of nouns in this toy corpus, the sample size is $N = 9$. There is one cooccurrence of *young* and *gentleman* ($O_{11} = 1$), two items where *gentleman* is modified by another adjective

	* w_2	* $\neg w_2$			* gent.	* \neg gent.	
$w_1 *$	O_{11}	O_{12}	$= f_1$	young *	1	2	$= 3$
$\neg w_1 *$	O_{21}	O_{22}		\neg young *	2	4	
	$= f_2$	$= N$			$= 3$	$= 9$	

Fig. 58.6: Contingency table of observed frequencies for syntactic cooccurrence, with concrete example for the word pair (*young, gentleman*) and the data in Figure 58.3 (right panel)

	$w_2 \in S$	$w_2 \notin S$			$\text{over} \in S$	$\text{over} \notin S$	
$w_1 \in S$	O_{11}	O_{12}	$= f_1$	$\text{hat} \in S$	1	2	$= 3$
$w_1 \notin S$	O_{21}	O_{22}		$\text{hat} \notin S$	1	1	
	$= f_2$	$= N$			$= 2$	$= 5$	

Fig. 58.7: Contingency table of observed frequencies for textual cooccurrence, with concrete example for the word pair (*hat*, *over*) and the data in Figure 58.2 (right panel)

($O_{21} = 2$), two items where *young* modifies another noun ($O_{12} = 2$), and four items that contain neither the adjective *young* nor the noun *gentleman* ($O_{22} = 4$).

For textual cooccurrence, Figure 58.2 motivates the definition of cooccurrence items as instances of textual units. In this example, each item corresponds to a sentence of the corpus. The sentence is classified into the first row of the contingency table if it contains one or more instances of w_1 and into the second row otherwise; it is classified into the first column if it contains one or more instances of w_2 and into the second column otherwise (see Figure 58.7). Note that no distinction is made between single and multiple occurrence of w_1 or w_2 in the same sentence. Again, the first row and column sums correspond to the marginal frequencies f_1 and f_2 as defined in section 3.2. The right panel of Figure 58.7 shows a contingency table for the word pair (*hat*, *over*), based on the example in Figure 58.2. With five sentences in the toy corpus, sample size is $N = 5$. One of the sentences contains both *hat* and *over* ($O_{11} = 1$), two sentences contain *hat* but not *over* ($O_{12} = 2$), one sentence contains *over* but not *hat* ($O_{21} = 1$), and one sentence contains neither of the two words ($O_{22} = 1$).

	w_2	$\neg w_2$			roll	$\neg \text{roll}$	
$\text{near}(w_1)$	O_{11}	O_{12}	$\approx k \cdot f_1$	$\text{near}(\text{hat})$	2	18	$= 20$
$\neg \text{near}(w_1)$	O_{21}	O_{22}		$\neg \text{near}(\text{hat})$	1	87	
	$= f_2$	$= N - f_1$			$= 3$	$= 108$	

Fig. 58.8: Contingency table of observed frequencies for surface cooccurrence, with concrete example for *roll* as a collocate of the node *hat* according to Figure 58.1 (right panel)

The statistical interpretation of surface cooccurrence is less straightforward than for the other two types. The most sensible definition identifies cooccurrence items with the relevant word tokens in the corpus, but excluding instances of the node word w_1 , for which no meaningful cross-classification is possible. Each item, i.e. word token, is then classified into the first row of the contingency table if it cooccurs with the node word w_1 , i.e. if it falls into one of the collocational spans around the instances of w_1 ; it is classified into the second row otherwise. The item is classified into the first column of the table if it is an instance of the targeted collocate w_2 , and into the second column otherwise. The

procedure is illustrated in Figure 58.8, with a concrete example for the data of Figure 58.1 shown in the right panel. This toy corpus consists of 111 word tokens (excluding punctuation). Subtracting the three instances of the node word *hat*, we obtain a sample size of $N = 108$. Of the 108 cooccurrence items, 20 fall into the collocational spans around instances of *hat*, so that the first row sum is $R_1 = 20$. Two of these items are cooccurrences of *hat* and *roll* ($O_{11} = 2$), and the remaining 18 items are classified into the second cell ($O_{12} = 18$). The 88 items outside the collocational spans are classified analogously: there is one instance of the collocate *roll* ($O_{21} = 1$), and all other items are assigned to the last cell of the table ($O_{22} = 87$).

For syntactic and textual cooccurrence, the contingency tables can be calculated directly from frequency signatures (O, f_1, f_2, N) that have been obtained as described in sections 3.2. and 3.3., using the following transformation equalities:

$$\begin{array}{ll} O_{11} = O & O_{12} = f_1 - O \\ O_{21} = f_2 - O & O_{22} = N - f_1 - f_2 + O \end{array}$$

The use of frequency signatures in combination with the equalities above is usually the most practical and convenient implementation of contingency tables. Tables for surface cooccurrence cannot be simplified in the same way, and it is recommended to calculate them by the explicit cross-classification procedure explained above.

5.2. Selected measures

Statistical association measures assume that the set of cooccurrence items is a random sample from a large population (representing an extensional definition of language as the set of all utterances that have been or can be produced, cf. article 36) and attempt to draw inferences about this population. Like simple measures, they can be divided into the general groups of effect-size and significance measures.

Effect-size measures aim to quantify how strongly the words in a pair are attracted to each other, i.e. they measure statistical association between the cross-classifying factors in the contingency table. Liebetrau (1983) gives a comprehensive survey of such association coefficients and Evert (2004, 54–58) discusses their mathematical properties. Coefficients describe properties of a population without taking sampling variation into account. They can be used as association measures in a straightforward way if this fact is ignored and the observed frequencies are taken as direct estimates for the corresponding population probabilities. As a result, effect-size measures tend to be unreliable especially for low-frequency data.

MI is the most intuitive association coefficient, comparing observed cooccurrence frequency against the value expected under the null hypothesis of independence. The equation shown in Figure 58.4 is also meaningful as a statistical association measure, where it should more precisely be written $\log_2(O_{11}/E_{11})$. Two other association coefficients are the (logarithmic) *odds ratio* (Blaheta/Johnson 2001, 56) and the *Dice coefficient* (Smadja/McKeown/Hatzivassiloglou 1996), shown in Figure 58.9. The odds ratio measure satisfies both conventions for association scores, with a value of 0 corresponding to independence and high positive values indicating strong positive association. Its interpretation is less intuitive than that of MI, though, and it has rarely been applied to

$$\begin{aligned} \text{chi-squared} &= \sum_{ij} \frac{(O_{ij} - E_{ij})^2}{E_{ij}} & \text{chi-squared}_{\text{corr}} &= \frac{N(|O_{11}O_{22} - O_{12}O_{21}| - N/2)^2}{R_1R_2C_1C_2} \\ \text{log-likelihood} &= 2 \sum_{ij} O_{ij} \log \frac{O_{ij}}{E_{ij}} & \text{average-MI} &= \sum_{ij} O_{ij} \cdot \log_2 \frac{O_{ij}}{E_{ij}} \\ \text{Dice} &= \frac{2O_{11}}{R_1 + C_1} & \text{odds-ratio} &= \log \frac{(O_{11} + \frac{1}{2})(O_{22} + \frac{1}{2})}{(O_{12} + \frac{1}{2})(O_{21} + \frac{1}{2})} \end{aligned}$$

Fig. 58.9: Some widely used statistical association measures

collocations. The Dice coefficient does not adhere to the second convention, as it does not assume a well-defined value in the case of independence. It cannot be used to identify word pairs with strong negative association, but is well-suited for rigid combinations such as fixed multiword units (Smadja/McKeown/Hatzivassiloglou 1996; Dias/Guilloré/Lopes 1999).

Statistical significance measures are based on the same types of hypothesis tests as the simple measures in section 4.2., viz. chi-squared tests (as a generalisation of z-scores) and likelihood-ratio tests. Unsurprisingly, there is no counterpart for t-score, which was based on an inappropriate test and hence cannot be translated into a rigorous statistical measure. The *chi-squared* measure adds up squared z-scores for all cells of the contingency table (\sum_{ij} indicates summation over all four cells, i. e. over indices $ij = 11, 12, 21, 22$). The normal approximation implicit in the z-scores becomes inaccurate if any of the expected frequencies E_{ij} are small, and chi-squared exhibits a low-frequency bias similar to the z-score measure. A better approximation is obtained by applying *Yates' continuity correction* (cf. DeGroot/Schervish 2002, section 5.8.). The continuity-corrected version is often written in the compact form shown as $\text{chi-squared}_{\text{corr}}$ in Figure 58.9, without explicit reference to expected frequencies E_{ij} . Chi-squared is a two-sided measure because the squared values are always positive. It can be transformed into a one-sided measure using the general procedure introduced in section 4.2. Chi-squared is often abbreviated X^2 , the symbol used for the chi-squared test statistic in mathematical statistics.

The *log-likelihood* measure (Dunning 1993) is a straightforward extension of simple-II, replacing the term $O - E$ by a summation over the remaining three cells of the contingency table. It is a two-sided measure and is sometimes abbreviated G^2 in analogy to X^2 . Interestingly, the association scores of log-likelihood, **simple-II** and chi-squared are all interpreted against the same scale, a χ^2 distribution (cf. section 4.2.). Mathematicians generally agree that the most appropriate significance test for contingency tables is *Fisher's exact test* (Agresti 2002, 91–93), which was put forward by Pedersen (1996) as an alternative to the log-likelihood measure. Unlike chi-squared and likelihood-ratio tests, this exact test does not rely on approximations that may be invalid for low-frequency data. Fisher's test can be applied as a one-sided or two-sided measure and provides a useful reference point for the discussion of other significance measures. However, it is computationally expensive and a sophisticated implementation is necessary to avoid numerical instabilities (Evert 2004, 93). Section 6.1. shows that log-likelihood provides an excellent approximation to association scores computed by Fisher's test, so there is little reason to use the complicated and technically demanding exact test. The informa-

tion-theoretic measure *average-MI* is identical to log-likelihood (except for a constant factor) and need not be discussed further here.

Note that simple association measures can also be computed from the full contingency tables, replacing O by O_{11} and E by E_{11} in the equations given in Figure 58.4. This shows clearly that many simple measures can be understood as a simplified version (or approximation) of a corresponding statistical measure. A more comprehensive list of association measures with further explanations can be found in Evert (2004, section 3) and online at:

<http://www.collocations.de/AM/>

Both resources describe simple as well as statistical association measures, using the notation for contingency tables introduced in this section and summarised in Figure 58.5.

6. Finding the right measure

The twelve equations in Figures 58.4 and 58.9 represent just a small selection of the many association measures that have been suggested and used over the years. Evert (2004) discusses more than 30 different measures, Pecina (2005) lists 57 measures, and new measures and variants are constantly being invented. While some measures have been established as de facto standards, e.g. log-likelihood in computational linguistics, t-score and MI in computational lexicography, there is no ideal association measure for all purposes. Different measures highlight different aspects of collocativity and will hence be more or less appropriate for different tasks: the n-best lists in Tables 58.2 and 58.3 are a case in point. The goal of this section is to help researchers choose a suitable association measure (or set of measures) for their study. While the primary focus is on understanding the characteristic properties of the measures presented in this article and the differences between them, the methods introduced below can also be applied to other association measures, allowing researchers to make an informed choice from the full range of options.

6.1. Mathematical arguments

Theoretical discussions of association measures are usually concerned with their mathematical derivation: the assumptions of the underlying model, the theoretical quantity to be measured, the validity and accuracy of the procedures used (especially if approximations are involved), and general mathematical properties of the measures (such as a bias towards low- or high-frequency word pairs). A first step in such discussions is to collect association measures with the same theoretical basis into groups. Measures within each group can often be compared directly with respect to their mathematical properties (since ideally they should measure the same theoretical quantity and hence lead to the same results), while different groups can only be compared at a general and rather philosophical level (does it make more sense to measure effect size or significance of association?).

As has already been mentioned in sections 4 and 5, the association measures introduced in this article fall into two major groups: *effect-size measures* (MI, Dice, odds ratio) and *significance measures* (z-score, t-score, simple-ll, chi-squared, log-likelihood). The choice between these two groups is largely a philosophical issue: one cannot be considered “better” than the other. Instead, they highlight different aspects of collocativity and are plagued by different types of mathematical problems.

Significance measures are particularly amenable to mathematical discussions, since in principle they attempt to measure the same theoretical quantity: the amount of evidence provided by a sample against the null hypothesis of independence. Moreover, chi-squared, log-likelihood and simple-ll use the same scale (the χ^2_1 distribution), so that their scores are immediately comparable. While z-score and t-score use a scale based on the normal distribution, their scores can easily be transformed to the χ^2_1 scale. The long-standing debate in mathematical statistics over appropriate significance tests for contingency tables has not completely been resolved yet (see Yates 1984), but most researchers consider Fisher’s exact test to be the most sensible and accurate measure of significance (Yates 1984, 446). We will therefore use it as a reference point for the comparison of association measures in the significance group. Fisher’s test calculates so-called p-values (cf. article 36), which are also transformed to the χ^2_1 scale for the comparison. The scatterplots in Figure 58.10 compare association scores calculated by various significance measures with those of Fisher’s exact test, using a synthetic data set in which cooccurrence and marginal frequencies have been varied systematically. The log-likelihood measure (G^2) and to some extent also simple-ll (G^2_{simple}) give an excellent approximation to Fisher’s test, as all data points are close to the diagonal. Chi-squared and z-score overestimate significance drastically (points far above diagonal), while t-score underestimates significance to a similar degree (points far below diagonal).

For effect-size measures, there is no well-defined theoretical quantity that would allow a direct comparison of their scores (e. g. with scatterplots as in Figure 58.10). Numerous coefficients have been suggested as measures of association strength in the population, but statisticians do not agree on a theoretically satisfactory choice (see e. g. Liebetrau 1983). A common mathematical property of effect-size measures is the use of direct estimates that do not take sampling variation into account. As a result, association

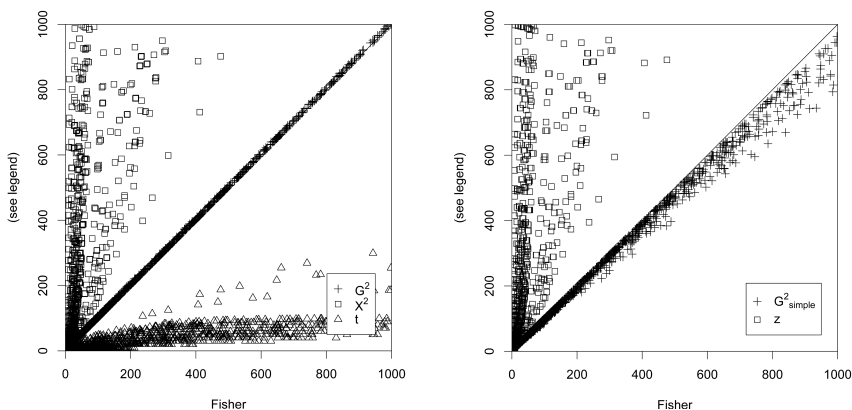


Fig. 58.10: Direct comparison of association scores on a synthetic data set, using Fisher’s exact test as a reference point (scores are transformed to χ^2_1 scale)

scores tend to become unreliable for low-frequency data. This effect is particularly severe for MI, odds ratio and similar measures that compare observed and expected frequency, since $E_{11} \ll 1$ for many low-frequency word pairs. Extending effect-size measures with a correction for sampling variation is a current topic of research and is expected to bridge the gap between the effect-size and significance groups (see section 7.1.).

It should be emphasised that despite their mathematical shortcomings, measures such as chi-squared and t-score may have linguistic merits that justify their use as heuristic measures for collocation identification. While clearly not satisfactory as measures of significance, they must not completely be excluded from the following discussion, which focuses on empirical and intuitive properties of association measures.

6.2. Collocations and multiword extraction

In those cases where mathematical theory does not help us choose between association measures, we can study their empirical properties independent of the underlying statistical reasoning. In this section, we specifically address empirical *linguistic* properties, i. e. we ask what kinds of word pairs are identified as collocations by the different association measures. A simple approach is to look at n-best lists as shown in Tables 58.2 and 58.3, which give a good impression of the different linguistic aspects of collocativity that the association measures capture. For instance, Table 58.2 indicates that simple-II is a useful measure for identifying typical and intuitively plausible collocates of a node word. Without a frequency threshold, MI brings up highly specialised terms (**-record bucket*), but also many obviously accidental cooccurrences (such as *dippermouth* or *Dok*). A more thorough and systematic study along these lines has been carried out by Stubbs (1995).

More precise empirical statements than such impressionistic case studies can be made if there is a well-defined goal or application for the identified collocations. A particularly profitable setting is the use of association scores for multiword extraction, where the goal is usually to identify a particular subtype of multiword expressions, e. g. compounds (Schone/Jurafsky 2001), technical terminology (Daille 1994) or lexical collocations (Krenn 2000). Evert/Krenn (2001, 2005) suggest an evaluation methodology for such tasks that allows fine-grained quantitative comparisons between a large number of association measures. The evaluation follows the standard procedure for semi-automatic multiword extraction, where recurrent word pairs are obtained from a corpus, optionally filtered by frequency or other criteria, and ranked according to a selected association measure. Since there are no meaningful absolute thresholds for association scores (cf. section 2.1.), it is standard practice to select an n-best list of the 500, 1000 or 2000 highest-ranking collocations as candidate multiword expressions. The candidates are then validated by an expert, e. g. a professional lexicographer or terminologist.

In the evaluation setting, candidates in the n-best list are manually annotated as *true positives* (i. e. multiword expressions of the desired type) and *false positives*. These annotations are used to calculate the *precision* of the n-best list, i. e. the proportion of true positives among the n multiword candidates, and sometimes also *recall*, i. e. how many of all suitable multiword expressions that could have been extracted from the corpus are actually found in the n-best list. The precision values of different association measures can then be compared: the higher the precision of a measure, the better it is

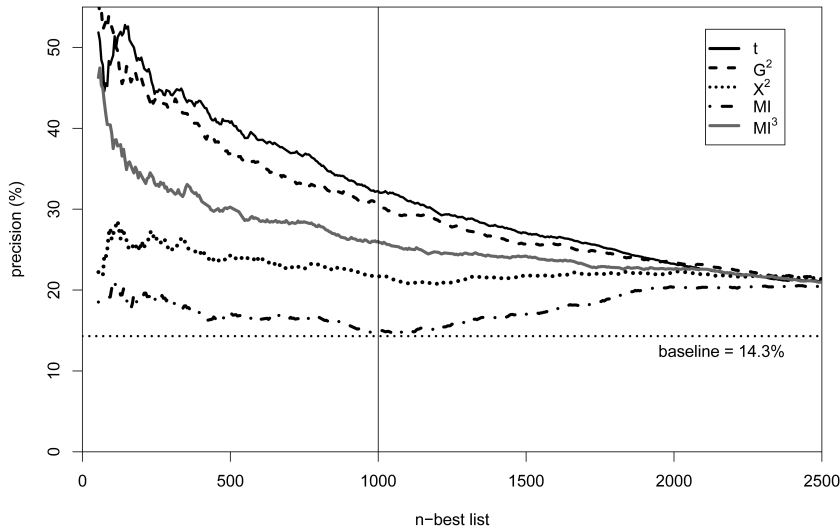


Fig. 58.11: Comparative evaluation of the association measures t -score (t), log-likelihood (G^2), chi-squared (X^2), MI and MI^3 on the data set of Krenn (2000)

suiting for identifying the relevant type of multiword expressions. Such evaluation experiments could be used, e.g., to confirm our impression that MI reliably identifies multiword proper names among adjacent bigrams (Table 58.3).

Instead of large and confusing tables listing precision values for various association measures and n -best lists, evaluation results can be presented in a more intuitive graphical form as *precision plots*. Figure 58.11 illustrates this evaluation methodology for the data set of Krenn (2000), who uses PP-verb cooccurrences from an 8-million-word subset of the German *Frankfurter Rundschau* newspaper corpus to identify lexical collocations between prepositional phrases and verbs (including support verb constructions and figurative expressions). The lines in Figure 58.11 summarise the precision values of five different association measures for arbitrary n -best lists. The precision for a particular n -best list can easily be read off from the graph, as indicated by the thin vertical line for $n = 1000$: the solid line at the top shows that t -score achieves a precision of approx. 32% on the 1000-best list, while log-likelihood (the dashed line below) achieves only 30.5%. The precision of chi-squared (dotted line) is much lower at 21.5%. Looking at the full lines, we see that log-likelihood performs much better than chi-squared for all n -best lists, as predicted by the mathematical discussion in section 6.1. Despite the frequency threshold, MI performs worse than all other measures and is close to the *baseline precision* (dotted horizontal line) corresponding to a random selection of candidates among all recurrent word pairs. Evaluation results always have to be interpreted in comparison to the baseline, and an association measure can only be considered useful if it achieves substantially better precision. The most striking result is that t -score outperforms all other measures, despite its mathematical shortcomings. This illustrates the limitations of a purely theoretical discussion: empirically, t -score is the best indicator for lexical PP-verb collocations among all association measures.

6.3. An intuitive geometrical model

In the previous section, we have looked at “linguistic” properties of association measures, viz. how accurately they can identify a particular type of multiword expressions or one of the other linguistic phenomena behind collocativity (see section 2.2.). If we take a pre-theoretic view of collocations as an observable property of language, though, the purpose of association scores is to measure this property in an appropriate way, not to match theoretical linguistic concepts. In this context, evaluation studies that depend on a theoretical or intuitive definition of true positives seem less appropriate. Instead, our goal should be to understand which quantitative aspects of collocativity each association measure singles out: we are interested in empirical mathematical properties of the measures.

Evert (2004, section 3.4.) proposes a geometric visualisation technique in order to reach the desired intuitive understanding of association measures. This technique works well for simple measures that require only two real numbers, O and E , to calculate an association score for a given word pair. Interpreting the numbers (O, E) as two-dimensional coordinates, we can thus represent each word pair in a data set by a point in the real Euclidean plane. The left panel of Figure 58.12 illustrates this “point cloud” view for adjacent bigrams in the Brown corpus. The data point representing the bigram *New York* (with $O = 303$ and $E \approx 0.54$) is marked with a circle. Its expected frequency $E \approx 0.5$ can be read off the x-axis, and its observed frequency $O = 303$ off the y-axis, as indicated by the thin horizontal and vertical lines. Note that both axes are on logarithmic scales in order to accommodate the wide range of observed and expected frequencies found in a typical data set. The frequency threshold $f \geq 10$ applied to the data set is clearly visible in the graph.

Association scores are usually compared against a cutoff threshold, whose value is either established in advance (in a threshold approach) or determined interactively (for n-best lists). In terms of the geometric model, the point cloud representing a data set is divided into *accepted* and *rejected* points by such a cutoff threshold. For any given association measure and cutoff threshold, this decision only depends on the coordinates of

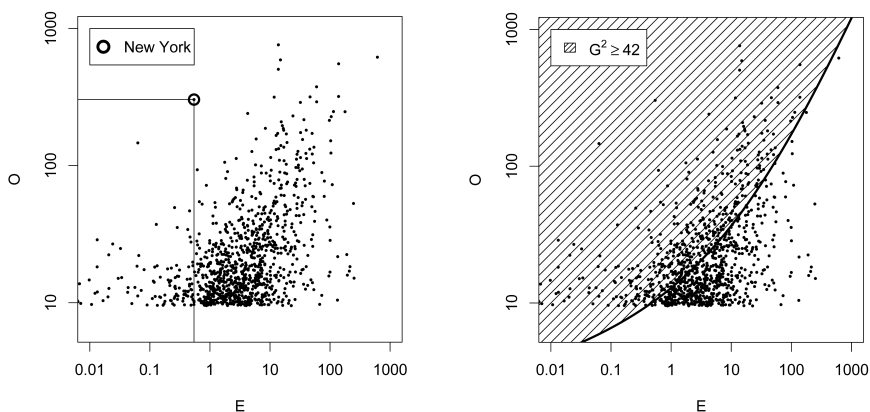


Fig. 58.12: Geometric visualisation of cooccurrence frequency data (left panel) and an acceptance region of the simple-II association measure (right panel)

a point in the Euclidean plane, not on the word pair represented by the point. It is therefore possible to determine for any hypothetical point in the plane whether it would be accepted or rejected, i.e. whether the association score would be higher than the threshold or not. The right panel of Figure 58.12 shows an illustration for the simple-II measure and a cutoff threshold of 42. Any data point in the shaded region will be assigned a score $G^2 \geq 42$, and any point outside the region a score $G^2 < 42$.

It can be shown that for most association measures the set of accepted hypothetical points forms a simple connected *acceptance region*. The region is bounded below by a single increasing line referred to as a *contour* of the association measure. All points on a contour line have the same association score according to this measure; in our example, a simple-II score of 42. Every simple association measure is uniquely characterised by its contour lines for different threshold values. We can thus visualise and compare measures in the form of contour plots as shown in Figure 58.13. Each panel overlays the contour plots of two different association measures. Comparing the shapes of the contour lines, we can identify the characteristic mathematical properties of the measures and understand the differences between them. Reading contour plots takes some practice: keep in mind that contours connect points with the same association scores, just as the contour lines of a topographic map connect points of the same elevation.

MI only considers the ratio between O and E , even for very low observed frequency O . Hence its dashed contours in Figure 58.13 are straight lines. These straight lines of constant ratio O/E also provide a grid for the interpretation of other contour plots. A significance measure such as simple-II (left panel) is sensitive to the smaller amount of evidence provided by low-frequency data. Therefore, a higher ratio between O and E is required to achieve the same score, and the contour lines have a left curvature. There is a single straight contour line, which marks the null hypothesis of independence ($O = E$) and coincides with the corresponding contour line of MI. Contours for positive association are located above and to the left of the independence line. Contours for negative association show a right curvature and are located below and to the right of the independence line.

The centre panel of Figure 58.13 shows a contour plot for the t-score measure. Again, independence is marked by a straight line that coincides with the MI contour. For positive association, the t-score contours have a left curvature similar to simple-II, but much more pronounced. For very small expected frequencies, they flatten out to horizontal lines, creating an implicit frequency threshold effect. We may speculate that this implicit

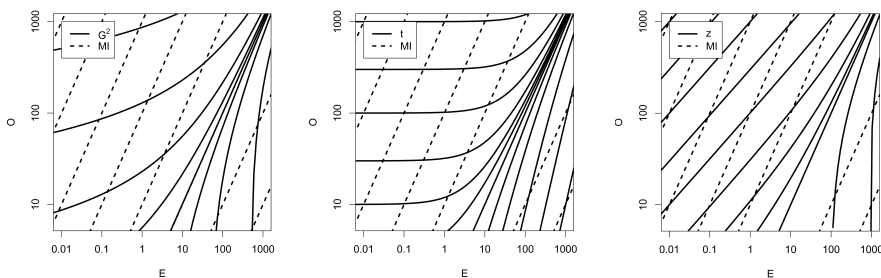


Fig. 58.13: Intuitive comparison of simple association measures represented by contour plots. The three panels compare simple-II (G^2 , left panel), t-score (centre panel) and z-score (right panel) against MI (dashed lines)

threshold is responsible for the good performance of t-score in some evaluation studies, especially if low-frequency word pairs are not discarded in advance. Interestingly, the contour lines for negative association are nearly parallel and do not seem to take random variation into account, in contrast to simple-II.

Finally, the right panel shows contour lines for z-score. Despite its mathematical background as a significance measure, z-score fails to discount low-frequency data. The contour lines for positive association are nearly parallel, although their slope is less steep than for MI. Thus, even data points with low observed frequency O can easily achieve high association scores, explaining the low-frequency bias of z-score that has been noted repeatedly. Interestingly, z-score seems to work well as a measure of significance for negative association, where its contour lines are very similar to those of simple-II.

The visualisation technique presented in this section can be extended to statistical association measures, but the geometric interpretation is more difficult and requires three-dimensional plots. See Evert (2004, section 3.4.) for details and sample plots.

7. Summary and conclusion

In this article, we have been concerned with the empirical Firthian notion of *collocations* as observations on the combinatorics of words in a language, which have to be distinguished clearly from lexicalised *multiword expressions* as pre-fabricated units, and in particular from *lexical collocations*, a subtype of multiword expressions. From the perspective of theoretical linguistics, collocations are often understood as an *epiphenomenon*, the surface reflections of compounds, idioms, lexical collocations and other types of multiword expressions, selectional preferences, semantic restrictions, cultural stereotypes, and to a considerable extent also conceptual knowledge (“facts of life”).

Introduced as an intuitively appealing, but fuzzy and pre-theoretical notion by Firth (1957), collocativity can be operationalised in terms of *cooccurrence* frequencies and quantified by mathematical *association measures*. High association scores indicate strong attraction between two words, but there is no standard scale of measurement to draw a clear distinction between collocations and non-collocations. Association measures and collocations have many uses, ranging from technical applications in computational linguistics to lexicographic and linguistic studies, where they provide descriptive generalisations about word usage. Collocations are closely related to lexicalised multiword expressions, and association measures are central to the task of automatic multiword extraction from corpora.

In order to identify and score collocations from a given corpus, the following steps have to be performed: (1) Choose an appropriate *type of cooccurrence* (surface, textual or syntactic). (2) Determine *frequency signatures* (i.e. cooccurrence frequency f and the marginal frequencies f_1 and f_2 in the corpus) for all relevant word pairs (w_1, w_2) as described in section 3 (Figures 58.1, 58.2 and 58.3 serve as a reminder), as well as sample size N . (3) Filter the cooccurrence data set by applying a *frequency threshold*. Theoretical considerations suggest a minimal threshold of $f \geq 3$ or $f \geq 5$, but higher thresholds often lead to even better results in practice. (4) Calculate the *expected frequencies* of the word pairs, using the general equation $E = f_1 f_2 / N$ for textual and syntactic cooccurrence, and the approximation $E = k f_1 f_2 / N$ for surface cooccurrence, where k is the total span size.

(5) Apply one of the *simple association measures* shown in Figure 58.4, or produce multiple tables according to different measures. Recall that the cooccurrence frequency f is denoted by O (for *observed frequency*) in these equations. (5) If collocations are treated as units, *rank* the word pairs by association score, or select a threshold to distinguish between collocations and non-collocations (or “strong” and “weak” collocations). In the node-collocate view, collocates w_2 are ranked separately for each node word w_1 .

If the data include word pairs with highly skewed marginal frequencies and you suspect that this may have distorted the results of the collocation analysis, you may want to apply *statistical association measures* instead of the simple measures. In order to do so, you have to compute a full 2×2 contingency table for each word pair, as well as a corresponding table of expected frequencies (see Figure 58.5). The precise calculation procedure depends on the type of cooccurrence and is detailed in section 5.1. (Figures 58.6, 58.7 and 58.8 serve as quick reminders). Then, one or more of the statistical measures in Figure 58.9 can be applied. Many further measures are found in Evert (2004) and online at <http://www.collocations.de/AM/> (both resources use the same notation as in this article).

The resulting set or ranking of collocations depends on many parameters, including the size and composition of the corpus, pre-processing (such as lemmatisation), application of frequency thresholds, the definition of cooccurrence used, and the choice of association measure. It is up to the researcher to find a suitable and meaningful combination of parameters, or to draw on results from multiple parameter settings in order to highlight different aspects of collocativity. While a particular type of cooccurrence is often dictated by the theoretical background of a study or practical restrictions (e.g., syntactic cooccurrence requires sufficiently accurate software for automatic syntactic analysis, or a pre-parsed corpus), other parameter values are more difficult to choose (e.g. span size for surface cooccurrence, or the frequency threshold).

A crucial step, of course, is to select one of well over 50 different association measures that are currently available (or to invent yet another measure). At this point, no definitive recommendation can be made. It is perhaps better to apply several measures with well-understood and distinct properties than attempt to find a single optimal choice. In any case, a thorough understanding of the characteristic properties of association measures and the differences between them is essential for a meaningful interpretation of the extracted collocations and their rankings. In section 6, various theoretical and empirical techniques have been introduced for this purpose, and the properties of several widely used measures have been discussed.

7.1. Open questions and extensions

The goal of this article was to present the current state of the art with regard to collocations and association measures. The focus has therefore been on established results rather than unsolved problems, open research questions, or extensions beyond simple word pairs. The following paragraphs give an overview of important topics of current research.

Like all statistical approaches in corpus linguistics, association measures suffer from the fact that the assumptions of their underlying statistical models are usually not met

by corpus data. In addition to the general question whether any finite corpus can be representative of a language (which is a precondition for the validity of statistical generalisations), *non-randomness* of corpus frequency data is a particularly serious problem for all statistical models based on random samples. A thorough discussion of this problem and possible solutions can be found in article 36 and in Evert (2006).

In addition to these common issues, cooccurrence data pose two specific difficulties. First, the null hypothesis of independence is extremely unrealistic. Words are never combined at random in natural language, being subject to a variety of syntactic, semantic and lexical restrictions. For a large corpus, even a small deviation from the null hypothesis may lead to highly significant rejection and inflated association scores calculated by significance measures. Effect-size measures are also subject to this problem and will produce inflated scores, e. g. for two rare words that always occur near each other (such as *déjà* and *vu*). A possible solution would be to specify a more realistic null hypothesis that takes some of the restrictions on word combinatorics into account, but research along these lines is still at a very early stage.

Second, word frequency distributions are highly skewed, with few very frequent types and a large number of extremely rare types. This property of natural language, often referred to as *Zipf's law* (see articles 37 and 41), is even more pronounced for cooccurrence data. In combination with the quantisation of observed frequencies (it is impossible to observe $O = 0.7$ cooccurrences), Zipf's law invalidates statistical corrections for sampling variation to the extent that accidental cooccurrences between low-frequency words may achieve very high association scores. An extensive study of this effect has resulted in the recommendation to apply a frequency threshold of $f \geq 5$ in order to weed out potentially spurious collocations (Evert 2004, chapter 4). Non-randomness effects may exacerbate the situation and necessitate even higher thresholds. Current research based on more sophisticated models of Zipfian frequency distributions aims to develop better correction techniques that are less drastic than a simple frequency threshold.

Intuitively, "mutual expectancies" often hold between more than two words. This is particularly obvious in the case of multiword expressions: *kick ... bucket* is always accompanied by the definite article *the*, *humble pie* usually occurs with *eat*, and the bigram *New York* is often followed by *City*. Applying association measures to word pairs will only bring up fragments of such larger collocations, and the missing pieces have to be filled in from the intuition of a linguist or lexicographer. It is therefore desirable to develop suitable measures for word triples and larger n -tuples. First attempts to formulate such measures are straightforward generalisations of the equations of MI (Lin 1998), log-likelihood (Zinsmeister/Heid 2003), or the Dice coefficient (da Silva/Lopes 1999). Obviously, a deep understanding of the mathematical properties of association measures for word pairs as well as their shortcomings is essential for a successful extension.

With the extension to n -word collocations, regular patterns become more noticeable: in addition to the well-known collocation *carry emotional baggage*, we also find *carry cultural, historical, ideological, intellectual, political, ... baggage* (some of them even more frequent than *emotional baggage*). This evidence suggests a productive *collocational pattern* of the form *carry Adj baggage*, with additional semantic restrictions on the adjective. Many instances of such patterns are too rare to be identified in corpora by statistical means, but would intuitively be considered as collocations by human speakers (think of *carry phraseological baggage*, for instance). There has been little systematic research on the productivity of collocations so far, notable exceptions being Lüdeling/Bosch (2003) and Stevenson/Fazly/North (2004).

Many collocations are intuitively felt to be *asymmetric*. For instance, in the bigram *the Iliad*, *the* is a more important collocate for *Iliad* than *Iliad* is for *the*. In the terminology of Kjellmer (1991), the bigram is left-predictive, but not right-predictive. Although such asymmetries are often reflected in skewed marginal frequencies (the collocation being more important for the less frequent word), hardly any of the known association measures make use of this information. Preliminary research suggests that measures of *directed association* could be based on the ratios O/f_1 and O/f_2 (as estimators for the conditional probability that w_1 is accompanied by w_2 and vice versa), or could be formulated by putting the association score of a word pair (w_1, w_2) in relation to the scores of all collocates of w_1 and w_2 , respectively (Michelbacher/Evert/Schütze 2007).

Although many association measures are available, there is still room for improvement and it would be desirable to develop measures with novel properties. Most existing measures fall into one of two major groups, viz. effect-size and significance measures. Both groups have their strengths and weaknesses: effect-size measures do not correct for sampling variation, while significance measures are biased towards high-frequency word pairs with small effect sizes (which tend to be uninteresting from a linguistic point of view). New association measures might be able to combine aspects of effect-size and significance measures, striking a balance between the low-frequency bias of the former and the high-frequency bias of the latter. First steps in this direction are summarised by Evert (2004, section 3.1.8.), but have not led to satisfactory results yet.

7.2. Further reading

Evert (2004) gives a more detailed account of statistical models for association in contingency tables and their limitations, together with a comprehensive inventory of association measures and methods for the comparison and evaluation of different measures. An online version of the inventory can be found at <http://www.collocations.de/AM/>. Contingency tables and the statistical tests that form the basis of many association measures are explained in standard textbooks on mathematical statistics (e.g. DeGroot/Schervish 2002). Advanced books (e.g. Agresti 2002) introduce more sophisticated models for the analysis of contingency tables. Although these models have not found widespread use as association measures yet, they may become important for the development of novel measures and their extension beyond simple word pairs.

Bartsch (2004) offers an insightful theoretical discussion of collocations and their properties, as well as an excellent overview of the various empirical and phraseological definitions of the term. Exemplary proponents of the two views are Sinclair (1991) and Sinclair et al. (2004) on the empirical side, and standard textbooks (e.g. Burger/Buhofer/Sialm 1982) for the phraseological view. Current research on collocations and multiword expressions is collected in the proceedings of ACL Workshops on Multiword Expressions (Daille/Williams 2001; Levin/Tokunaga/Lenci 2003; Tanaka et al. 2004; Villada Moirón et al. 2006; Grégoire/Evert/Kim 2007) and in Grossmann/Tutin (2003).

Relevant articles in this volume are article 24 (on word segmentation and part-of-speech tagging), article 25 (on lemmatisation), article 26 (on word sense disambiguation) and article 28 (on automatic syntactic annotation), as well as article 10 (on text corpora). Article 36 is a general introduction to the statistical analysis of corpus frequency data,

including most of the techniques on which association measures are based. Important applications of collocations can be found in the articles on computational lexicography (article 8) and word meaning (article 45).

We have followed a traditional view of collocations as simple word pairs here, but association measures and related techniques can equally well be applied to cooccurrences of other linguistic units (e.g. lexical items and syntactic constructions in article 43).

8. Literature

- Agresti, Alan (2002), *Categorical Data Analysis*, 2nd edition. Hoboken: John Wiley & Sons.
- Aston, Guy/Burnard, Lou (1998), *The BNC Handbook: Exploring the British National Corpus with SARA*. Edinburgh: Edinburgh University Press. See also the BNC homepage at <http://www.natcorp.ox.ac.uk/>.
- Bartsch, Sabine (2004), *Structural and Functional Properties of Collocations in English*. Tübingen: Narr.
- Berry-Rogghe, Godelieve L. M. (1973), The Computation of Collocations and their Relevance to Lexical Studies. In: Aitken, Adam J./Bailey, Richard W./Hamilton-Smith, Neil (eds.), *The Computer and Literary Studies*. Edinburgh: Edinburgh University Press, 103–112.
- Blaheta, Don/Johnson, Mark (2001), Unsupervised Learning of Multi-word Verbs. In: *Proceedings of the ACL Workshop on Collocations*. Toulouse, France, 54–60.
- Burger, Harald/Buhofer, Annelies/Sialm, Ambros (1982), *Handbuch der Phraseologie*. Berlin etc.: Walter de Gruyter.
- Choueka, Yaacov (1988), Looking for Needles in a Haystack. In: *Proceedings of RIAO '88*. Cambridge, MA, 609–623.
- Church, Kenneth W./Hanks, Patrick (1990), Word Association Norms, Mutual Information, and Lexicography. In: *Computational Linguistics* 16(1), 22–29.
- Church, Kenneth/Gale, William A./Hanks, Patrick/Hindle, Donald (1991), Using Statistics in Lexical Analysis. In: Zernick, Uri (ed.) *Lexical Acquisition: Using On-line Resources to Build a Lexicon*. Hillsdale, NY: Lawrence Erlbaum, 115–164.
- da Silva, Joaquim Ferreira/Lopes, Gabriel Pereira (1999), A Local Maxima Method and a Fair Dispersion Normalization for Extracting Multi-word Units from Corpora. In: *6th Meeting on the Mathematics of Language*. Orlando, FL, 369–381.
- Daille, Béatrice (1994), *Approche mixte pour l'extraction automatique de terminologie: statistiques lexicales et filtres linguistiques*. PhD thesis, Université Paris 7.
- Daille, Béatrice/Williams, Geoffrey (eds.) (2001), *Proceedings of the 2001 ACL Workshop on Collocation*. Toulouse, France.
- DeGroot, Morris H./Schervish, Mark J. (2002), *Probability and Statistics*, 3rd edition. Boston: Addison Wesley.
- Dennis, Sally F. (1965), The Construction of a Thesaurus Automatically from a Sample of Text. In: Stevens, Mary E./Giuliano, Vincent E./Heilprin, Lawrence B. (eds.), *Proceedings of the Symposium on Statistical Association Methods for Mechanized Documentation*. (National Bureau of Standards Miscellaneous Publication 269.) Washington: National Bureau of Standards, 61–148.
- Dias, Gaël/Guilloré, Sylvie/Lopes, José G. P. (1999), Language Independent Automatic Acquisition of Rigid Multiword Units from Unrestricted Text Corpora. In: *Proceedings of Traitement Automatique des Langues Naturelles (TALN)*. Cargèse, Corsica, France, 333–338.
- Dunning, Ted E. (1993), Accurate Methods for the Statistics of Surprise and Coincidence. In: *Computational Linguistics* 19(1), 61–74.
- Evert, Stefan (2004), *The Statistics of Word Cooccurrences: Word Pairs and Collocations*. Dissertation, Institut für maschinelle Sprachverarbeitung, University of Stuttgart. Published in 2005, URN urn:nbn:de:bsz:93-opus-23714.

- Evert, Stefan (2006), How Random is a Corpus? The Library Metaphor. In: *Zeitschrift für Anglistik und Amerikanistik* 54(2), 177–190.
- Evert, Stefan/Kermes, Hannah (2003), Experiments on Candidate Data for Collocation Extraction. In: *Companion Volume to the Proceedings of the 10th Conference of the European Chapter of the Association for Computational Linguistics*. Budapest, Hungary, 83–86.
- Evert, Stefan/Krenn, Brigitte (2001), Methods for the Qualitative Evaluation of Lexical Association Measures. In: *Proceedings of the 39th Annual Meeting of the Association for Computational Linguistics*. Toulouse, France, 188–195.
- Evert, Stefan/Krenn, Brigitte (2005), Using Small Random Samples for the Manual Evaluation of Statistical Association Measures. In: *Computer Speech and Language* 19(4), 450–466.
- Firth, John Rupert (1957), A Synopsis of Linguistic Theory 1930–55. In: *Studies in Linguistic Analysis*. Oxford: The Philological Society, 1–32. Reprinted in Palmer 1968, 168–205.
- Gil, Alexandre/Dias, Gaël (2003), Using Masks, Suffix Array-based Data Structures and Multidimensional Arrays to Compute Positional Ngram Statistics from Corpora. In: *Proceedings of the ACL Workshop on Multiword Expressions*. Sapporo, Japan, 25–32.
- Goldman, Jean-Philippe/Nerima, Luka/Wehrli, Eric (2001), Collocation Extraction Using a Syntactic Parser. In: *Proceedings of the ACL Workshop on Collocations*. Toulouse, France, 61–66.
- Grégoire, Nicole/Evert, Stefan/Kim, Su Nam (eds.) (2007), *Proceedings of the Workshop on A Broader Perspective on Multiword Expressions*. Prague, Czech Republic.
- Grossmann, Francis/Tutin, Agnès (eds.) (2003), *Les collocations: Analyse et traitement*. Amsterdam: De Werelt.
- Hausmann, Franz Josef (1989), Le dictionnaire de collocations. In: Hausmann, Franz Josef/Reichmann, Otto/Wiegand, Herbert Ernst (eds.) *Wörterbücher, Dictionaries, Dictionnaires. Ein internationales Handbuch*. Berlin: Mouton de Gruyter, 1010–1019.
- Heid, Ulrich/Evert, Stefan/Docherty, Vincent/Worsch, Wolfgang/Wermke, Matthias (2000), A Data Collection for Semi-automatic Corpus-based Updating of Dictionaries. In: Heid, Ulrich/Evert, Stefan/Lehmann, Egbert/Rohrer, Christian (eds.), *Proceedings of the 9th EURALEX International Congress*. Stuttgart, Germany, 183–195.
- Kilgariff, Adam/Rychly, Pavel/Smrz, Pavel/Tugwell, David (2004), The Sketch Engine. In: *Proceedings of the 11th EURALEX International Congress*. Lorient, France, 105–116.
- Kjellmer, Göran (1991), A Mint of Phrases. In: Aijmer, Karin/Altenberg, Bengt (eds.), *English Corpus Linguistics*. London: Longman, 111–127.
- Krenn, Brigitte (2000), *The Usual Suspects: Data-oriented Models for the Identification and Representation of Lexical Collocations*. (Saarbrücken Dissertations in Computational Linguistics and Language Technology 7.) Saarbrücken: DFKI & Universität des Saarlandes.
- Lea, Diana (ed.) (2002), *Oxford Collocations Dictionary for Students of English*. Oxford etc.: Oxford University Press.
- Levin, Lori/Tokunaga, Takenobu/Lenci, Alessandro (eds.) (2003), *Proceedings of the ACL 2003 Workshop on Multiword Expressions: Analysis, Acquisition and Treatment*. Sapporo, Japan.
- Liebetrau, Albert M. (1983), *Measures of Association*. (Sage University Papers Series on Quantitative Applications in the Social Sciences 32.) Newbury Park: Sage.
- Lin, Dekang (1998), Extracting Collocations from Text Corpora. In: *Proceedings of the First Workshop on Computational Terminology*. Montreal, Canada, 57–63.
- Lüdeling, Anke/Bosch, Peter (2003), Identification of Productive Collocations. In: *Proceedings of the 8th International Symposium on Social Communication*. Santiago de Cuba, Cuba.
- Manning, Christopher D./Schütze, Hinrich (1999), *Foundations of Statistical Natural Language Processing*. Cambridge, MA: MIT Press.
- Michelbacher, Lukas/Evert, Stefan/Schütze, Hinrich (2007), Asymmetric Association Measures. In: *Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2007)*. Borovets, Bulgaria. Available at: <http://www.cogsci.uni-osnabrueck.de/~severt/PUB/MichelbacherEtc2007.pdf>.
- Palmer, Frank R. (ed.) (1968), *Selected Papers of J. R. Firth 1952–59*. London: Longmans.

- Pecina, Pavel (2005), An Extensive Empirical Study of Collocation Extraction Methods. In: *Proceedings of the ACL Student Research Workshop*. Ann Arbor, MI, 13–18.
- Pedersen, Ted (1996), Fishing for Exactness. In: *Proceedings of the South-Central SAS Users Group Conference*. Austin, TX, 188–200.
- Sag, Ivan A./Baldwin, Timothy/Bond, Francis/Copestake, Ann/Flickinger, Dan (2002), Multiword Expressions: A Pain in the Neck for NLP. In: *Proceedings of the Third International Conference on Intelligent Text Processing and Computational Linguistics (CICLing 2002)*. Mexico City, Mexico, 1–15.
- Sahlgren, Magnus (2006), *The Word Space Model: Using Distributional Analysis to Represent Syntagmatic and Paradigmatic Relations between Words in High-dimensional Vector Spaces*. PhD thesis, Department of Linguistics, Stockholm University.
- Schiehlen, Michael (2004), Annotation Strategies for Probabilistic Parsing in German. In: *Proceedings of COLING 2004*. Geneva, Switzerland, 390–396.
- Schone, Patrick/Jurafsky, Daniel (2001), Is Knowledge-free Induction of Multiword Unit Dictionary Headwords a Solved Problem? In: *Proceedings of the 2001 Conference on Empirical Methods in Natural Language Processing*. Pittsburgh, PA, 100–108.
- Schütze, Hinrich (1998), Automatic Word Sense Discrimination. In: *Computational Linguistics* 24(1), 97–123.
- Sinclair, John (1966), Beginning the Study of Lexis. In: Bazell, Charles E./Catford, John C./Halliday, Michael A. K./Robins, Robert H. (eds.), *In Memory of J. R. Firth*. London: Longmans, 410–430.
- Sinclair, John (1991), *Corpus, Concordance, Collocation*. Oxford: Oxford University Press.
- Sinclair, John (ed.) (1995), *Collins COBUILD English Dictionary*. New edition, completely revised. London: Harper Collins.
- Sinclair, John/Jones, Susan/Daley, Robert/Krishnamurthy, Ramesh (2004), *English Collocation Studies: The OSTI Report*. London etc.: Continuum Books. Originally written in 1970 (unpublished).
- Smadja, Frank (1993), Retrieving Collocations from Text: Xtract. In: *Computational Linguistics* 19(1), 143–177.
- Smadja, Frank/McKeown, Kathleen R./Hatzivassiloglou, Vasileios (1996), Translating Collocations for Bilingual Lexicons: A Statistical Approach. In: *Computational Linguistics* 22(1), 1–38.
- Stevenson, Suzanne/Fazly, Afsaneh/North, Ryan (2004), Statistical Measures of the Semi-productivity of Light Verb Constructions. In: *Proceedings of the ACL 2004 Workshop on Multiword Expressions: Integrating Processing*. Barcelona, Spain, 1–8.
- Stubbs, Michael (1995), Collocations and Semantic Profiles: On the Cause of the Trouble with Quantitative Studies. In: *Functions of Language* 2(1), 23–55.
- Tanaka, Takaaki/Villavicencio, Aline/Bond, Francis/Korhonen, Anna (eds.) (2004), *Proceedings of the Second ACL Workshop on Multiword Expressions: Integrating Processing*. Barcelona, Spain.
- Terra, Egidio/Clarke, Charles L. A. (2004), Fast Computation of Lexical Affinity Models. In: *Proceedings of COLING 2004*. Geneva, Switzerland, 1022–1028.
- Villada Moirón, Begoña/Villavicencio, Aline/McCarthy, Diana/Evert, Stefan/Stevenson, Suzanne (eds.) (2006), *Proceedings of the ACL Workshop on Multiword Expressions: Identifying and Exploiting Underlying Properties*. Sydney, Australia.
- Williams, Geoffrey (2003), Les collocations et l'école contextualiste britannique. In: Grossmann/Tutin 2003, 33–44.
- Yates, Frank (1984), Tests of Significance for 2 × 2 Contingency Tables. In: *Journal of the Royal Statistical Society, Series A* 147(3), 426–463.
- Zinsmeister, Heike/Heid, Ulrich (2003), Significant Triples: Adjective + Noun + Verb Combinations. In: *Proceedings of the 7th Conference on Computational Lexicography and Text Research (COMPLEX 2003)*. Budapest, Hungary, 92–101.

Stefan Evert, Osnabrück (Germany)