research directions, paying particular attention to platform governance as its impacts digital scholarship.

## Bibliography

**Chen, G., Posada, A., & Chan, L.** (2019). Vertical integration in academic publishing: implications for knowledge inequality." In Chan, L. and Mounier, P. (eds.), Connecting the Knowledge Commons—From Projects to Sustainable Infrastructure: The 22nd International Conference on Electronic Publishing – Revised Selected Papers. OpenEdition Press. <u>http://books.openedition.org/oep/9068</u>.

**Fitzpatrick, K**. (2018). Generous Thinking. Baltimore: Johns Hopkins University Press.

Liu, A. (2018). Toward critical infrastructure studies. NASSR: 1-22. <u>https://cistudies.org/wp-content/uploads/</u> Toward-Critical-Infrastructure-Studies.pdf.

**Pawlicka-Deger, U**. (2021). Infrastructuring digital humanities. Digital Scholarship in the Humanities, fqab086. <u>https://doi.org/10.1093/llc/fqab086</u>.

Plantin, J., Lagoze, C., Edwards, P. N., and Sandvig, C. (2018). Infrastructure studies meet platform studies in the age of Google and Facebook. New Media & Society, 20(1): 293-310.

van Dijck, J., Poell, T., and de Waal, M. (2018). The Platform Society: Public Values in a Connective World. New York: Oxford University Press.

### Measuring Keyness

#### **Evert, Stephanie**

stephanie.evert@fau.de FAU Erlangen-Nürnberg, Germany

## Keywords in corpus linguistics and DH

In corpus linguistics, the notion of **keywords** refers to words (and sometimes also multiword units, semantic categories or lexico-grammatical constructions) that "occur with unusual frequency in a given text" (Scott, 1997: 236) or a text collection, i.e. a corpus. Keywords are deemed to represent the characteristic vocabulary of the target text or corpus and thus have many applications in corpus linguistics, digital humanities and computational social science. They can capture the aboutness of a text (Scott, 1997), the terminology of a text genre or technical domain (Paquot and Bestgen, 2009), important aspects of literary style (Culpeper, 2009), linguistic and cultural differences (Oakes and Farrow, 2006), etc.; they give insight into historical perspectives (Fidler and Cvrček, 2015) and provide a basis for measuring the similarity of text collections (Rayson and Garside, 2000). Keywords are also an important starting point for corpus-based discourse analysis (Baker, 2006), where manually formed clusters of keywords represent central topics, actors, metaphors, and framings (e.g. McEnery et al., 2015). Since this process is guided from the outset by human understanding, it provides a more interpretable alternative to topic models in hermeneutic text analysis.

Keywords are usually operationalised in terms of a statistical frequency comparison between the **target corpus** and a **reference corpus**. Different research questions can be addressed depending on the particular constellation of target *T* and reference *R*, e.g. (i) *T* = a single text vs. *R* = a text collection ( $\rightarrow$  aboutness), (ii) *T* and *R* = collections of articles on the same topic in left-leaning and right-leaning newspapers ( $\rightarrow$  contrastive framings), or (iii) *T* = texts from a given domain or genre vs. *R* = a large general-language reference corpus ( $\rightarrow$  terminology).

Although keyword analysis is a well-established approach and has been implemented in many standard corpus-linguistic software tools such as WordSmith<sup>1</sup>, AntConc<sup>2</sup>, SketchEngine<sup>3</sup>, and CQPweb (Hardie, 2012), it is still unclear what the "right" way of measuring keyness is (see overview in Hardie, 2014). In this paper, I propose (i) a mathematically well-founded **best-practice technique** and (ii) introduce a **visual approach** for exploring the empirical properties of different keyness measures.

#### Keyness measures

Keyword analysis is operationalised as a comparison of relative frequencies: For each **candidate** word, its frequency  $f_1$  in a target corpus T of  $n_1$  tokens is compared to its frequency  $f_2$  in a reference corpus R of  $n_2$  tokens. The candidate set of m items typically includes words that only occur in the target corpus  $(f_2 = 0)$ .

A candidate is considered a ("positive") keyword if its relative frequency  $p_1 = f_1 / n_1$  in *T* is substantially higher than its relative frequency  $p_2 = f_2 / n_2$  in *R*. A large number of **keyness measures** have been proposed to quantify the comparison and thus provide a basis for a ranking of the candidates and/or cut-off thresholds. Three main groups of measures can be distinguished:

1. Measures based on **hypothesis tests** put the focus on establishing a statistically significant difference between  $p_1$  and  $p_2$ . The most widely-used measures are chi-

squared  $X^2$  and log-likelihood  $G^2$  (Dunning, 1993). These measures are biased towards high-frequency keywords, often including function words and other non-specific words.

- 2. Effect size measures instead focus on how many times more frequent a candidate is in *T* than in *R*. The most intuitive measure is relative risk  $r = p_1 / p_2$ , also known as LogRatio = log 2 *r* (Hardie, 2014). Some other effectsize measures are equivalent (%DIFF, Gabrielatos and Marchi, 2012) or closely related (odds ratio, Pojanapunya and Watson Todd, 2018) to LogRatio. These measures are biased towards very low-frequency keywords and are often combined with an additional significance filter (typically based on *G* <sup>2</sup>).
- 3. Various **heuristic** measures lack any statistical foundation. They are often particularly easy to compute such as SketchEngine's SimpleMaths (Kilgarriff, 2009), which also offers a user parameter to adjust its bias towards high-frequency or low-frequency keywords.

# Mathematical discussion and visualisation

Hypothesis-test measures are subject to the criticism raised more generally against p-value testing in corpus linguistics and other fields (e.g. Gries, 2005). In particular, they are biased towards high-frequency keywords irrespective of effect size, selecting candidates that are not very salient for the target corpus. When they are applied more reasonably as a significance filter, the problem of multiple testing is often ignored: a single analysis may carry out frequency comparisons for hundreds of thousands of candidates, resulting in large numbers of false positives at customary significance levels such as p < .001 (Gries, 2005; Hardie, 2014).

By contrast, effect-size measures such as LogRatio are biased towards low-frequency keywords because they completely ignore the statistical significance of the observed difference in relative frequency. Moreover, many of these measures are undefined for f = 0 and need special heuristics for this case; e.g. Hardie (2014) simply substitutes f = 0.5 without mathematical justification.

Traditionally, keyness measures are computed from cumulative token frequency counts for *T* and *R*. However, two recent studies have independently concluded that keywords based on document counts are more robust (Evert et al., 2018; Egbert and Biber, 2019).

Keyness measures can also be understood from a more intuitive angle by visualising them as **topographic maps**, which show the scores assigned to all possible combinations of frequencies  $f_1$  in T and  $f_2$  in R on a logarithmic scale (similar to the visualisation of collocations in Evert, 2004:

sec. 3.3). The examples in Fig. 1 reveal the respective frequency biases of  $G^2$  and LogRatio – which is hardly mitigated by an additional significance filter – in the top row (dark red colours indicate frequency profiles of highly-ranked keywords).



Visualisation of keyness measures as topographic maps for n = n = 100 M words. The bottom right panel highlights problems of an earlier version of LRC currently used by CQPweb.

### Best-practice recommendation

Conservative estimates based on statistical confidence intervals combine the advantages of hypothesis tests and effect-size measures into a single score. I therefore propose LRC, a conservative estimate of LogRatio, as a bestpractice keyness measure. LRC uses an exact conditional Poisson test (Fay, 2010: 55) to obtain reliable confidence intervals corrected for multiple testing. The full procedure for computing LRC scores is as follows:

- 1. Collect the frequency data  $f_1, f_2$  for each candidate and the sample sizes  $n_1, n_2$  of T and R. Wherever suitable, document frequencies should be preferred.
- Compute a two-sided Pearson-Clopper binomial confidence interval [π –, π +] for f 1 successes out of f 1 + f 2 trials, with Bonferroni-adjusted significance level α = 0.05 / m.
- 3. Convert the binomial proportions to  $[LRC -, LRC +] = [\log 2 (n 2 \pi / n 1 (1 \pi -)), \log 2 (n 2 \pi + / n 1 (1 \pi +))].$
- 4. If the test is not significant (LRC  $\le 0 \le LRC +$ ), set LRC = 0. Otherwise, set LRC = LRC if  $p_1 > p_2$  and LRC = LRC + if  $p_1 < p_2$ .

LRC has several **advantages** over other keyness measures: (i) it balances out the high-frequency bias of

hypothesis tests and the low-frequency bias of effect-size measures (cf. right panel of Fig. 2); (ii) unlike heuristics such as SimpleMaths it does this in a mathematically welljustified way; (iii) it can be applied to candidates with  $f_2 = 0$  without special precautions; (iv) it detects both positive ( $p_1 > p_2$ ) and negative ( $p_1 < p_2$ ) keywords; (v) it includes a reliable significance filter (LRC = 0) and does not require arbitrary frequency thresholds; (vi) robust and efficient implementations of the underlying binomial confidence intervals are available in standard statistical software packages, so very large candidate sets can easily be processed. The left panel of Fig. 2 shows that LRC overlaps well with established keyness measures, again indicating that it provides an excellent compromise.

A reference implementation of LRC is available at <u>https://osf.io/cy6mw/</u> together with a more detailed analysis. It is also included in version 0.6 of the *corpora* package for R. 4



Quantitative analysis of top-250 keyword lists for the data of Evert et al. (2018): overlap between four measures (left panel) and frequency distribution in the target corpus (right panel).

## Bibliography

**Baker, P.** (2006). *Using Corpora in Discourse Analysis*. London: Continuum Books.

**Culpeper, J.** (2009). Keyness: Words, parts-ofspeech and semantic categories in the character-talk of Shakespeare's Romeo and Juliet. *International Journal of Corpus Linguistics*, **14**(1): 29–59.

**Dunning, T. E.** (1993). Accurate methods for the statistics of surprise and coincidence. *Computational Linguistics*, **19**(1): 61–74.

**Egbert, J. and Biber, D.** (2019). Incorporating text dispersion into keyword analyses. *Corpora*, **14**(1): 77–104.

**Evert, S.** (2004). *The statistics of word cooccurrences: Word pairs and collocations*. Dissertation, Institut für maschinelle Sprachverarbeitung, University of Stuttgart, doi: . Evert, S., Dykes, N. and Peters, J. (2018). A quantitative evaluation of keyword measures for corpusbased discourse analysis. Presentation at the Corpora & Discourse International Conference (CAD 2018), Lancaster, UK, <u>https://www.stephanie-evert.de/PUB/ EvertEtc2018\_CAD\_slides.pdf</u>.

**Fay, M. P.** (2010). Two-sided exact tests and matching confidence intervals for discrete data. *The R Journal*, **2**(1): 53–58.

Fidler, M. and Cvrček, V. (2015). A data-driven analysis of reader viewpoints: Reconstructing the historical reader using keyword analysis. *Journal of Slavic Linguistics*, **23**(3): 197–239.

**Gabrielatos, C. and Marchi, A.** (2012). Keyness: Appropriate metrics and practical issues Presentation at the Corpora and Discourse Studies Conference (CADS 2012), Bologna, Italy, <u>https://www.researchgate.net/</u> publication/261708842\_Keyness\_Appropriate\_metrics\_ and\_practical\_issues.

**Gries, S. Th.** (2005). Null-hypothesis significance testing of word frequencies: A follow-up on Kilgarriff. *Corpus Linguistics and Linguistic Theory*, **1**(2): 277–94.

Hardie, A. (2012). CQPweb – combining power, flexibility and usability in a corpus analysis tool. *International Journal of Corpus Linguistics*, **17**(3): 380–409.

**Hardie**, **A.** (2014). A single statistical technique for keywords, lockwords, and collocations.

Kilgarriff, A. (2009). Simple maths for keywords. *Proceedings of the Corpus Linguistics 2009 Conference*. Liverpool, UK, <u>http://ucrel.lancs.ac.uk/publications/</u> <u>CL2009/</u>.

McEnery, T., McGlashan, M. and Love, R. (2015). Press and social media reaction to ideologically inspired murder: The case of Lee Rigby. *Discourse and Communication*, 9(2): 1–23, doi: 10.1177/1750481314568545.

**Oakes, M. P. and Farrow, M.** (2006). Use of the chisquared test to examine vocabulary differences in English language corpora representing seven different countries. *Literary and Linguistic Computing*, **22**(1): 85–99, doi: <u>10.1093/Ilc/fql044</u>.

**Paquot, M. and Bestgen, Y.** (2009). Distinctive words in academic writing: A comparison of three statistical tests for keyword extraction. In Jucker, A., Schreier, D. and Hundt, M. (eds), *Corpora: Pragmatics and Discourse. Papers from the 29th International Conference on English Language Research on Computerized Corpora*. Amsterdam: Rodopi, pp. 247–69, doi: 10.1163/9789042029101\_014.

**Pojanapunya, P. and Watson Todd, R.** (2018). Loglikelihood and odds ratio: Keyness statistics for different purposes of keyword analysis. *Corpus Linguistics and Linguistic Theory*, **14**(1): 133–67. **Rayson, P. and Garside, R.** (2000). Comparing corpora using frequency profiling. *Proceedings of the ACL Workshop on Comparing Corpora*. Hong Kong, pp. 1–6.

Scott, M. (1997). PC analysis of key words – and key key words. *System*, **25**(2): 233–45.

#### Notes

- 1. https://www.lexically.net/wordsmith/
- 2. https://www.laurenceanthony.net/software/antconc/
- 3. https://www.sketchengine.eu/
- 4. <u>https://cran.r-project.org/web/packages/corpora/</u>

## Historical Research meets Semantic Interoperability: The Documentation System SYNTHESIS and its Application in Art History Research

#### **Fafalios**, Pavlos

fafalios@ics.forth.gr Centre for Cultural Informatics, Institute of Computer Science, FORTH, Greece

## Introduction and motivation

Historical science is the field that describes, examines and questions a sequence of past events, and investigates patterns of cause and effect. Research in the field usually starts by first discovering, collecting, documenting and organizing historical sources, such as written documents or material artifacts. This often includes either the transcription (and then curation) of historical archival sources, like in Petrakis et al. (2020) for the case of Maritime History, or the detailed documentation of cultural artifacts and related evidence, like in Fafalios et al. (2021) for the case of Art History, with the latter being the focus of this presentation.

In this context, although computing in the field has developed enormously over the last years, data management problems still exist and are very varied. Common problems include: a) the difficulty for collaborative but controlled documentation by a large number of historians of different research groups; b) the lack of representation of the details from which the documented relations are inferred, important for the long-term validity of the research results; c) the difficulty to combine and integrate information extracted from multiple and diverse information sources; d) the difficulty of third parties to understand and re-use the documented data, resulting in the production of data with limited longevity.

## The SYNTHESIS system

In an effort to cope with the aforementioned problems, we present the SYNTHESIS documentation system and its use by a large number of historians in the context of a large European research project of Art History, called RICONTRANS (ERC Consolidator Grant, No 818791). SYNTHESIS is Web-based, multilingual, configurable (for use in other digital humanities fields), and utilizes XML technology, offering flexibility in terms of versioning, workflow management and data model extension. It focuses on semantic interoperability (Ouksel and Sheth, 1999), enabling the exchange of data among computer systems with unambiguous/shared meaning, and achieves this by making use of standards for data modelling and publication, in particular the formal ontology CIDOC-CRM (ISO 21127:2014) and the data model RDF (W3C Recommendation). The aim is the production of data with high value, longevity and long-term validity that can be (re)used beyond a particular research activity.

SYNTHESIS offers a wide range of functionalities including i) interlinking of the documented entities (forming a network of interrelated entities), ii) management of static and dynamic vocabularies, iii) linking to thesauri of terms, iv) connection with geolocation services (TGN, Geonames), v) map visualization for certain types of entities, vi) support of comparable historical time expressions (e.g., ca. 1920, early 16th century), vii) management of digital files (images, etc.), viii) transformation of the documented information to a knowledge base of Linked Data (Heath and Bizer, 2011).

VISUAL CULTURE PETY AND PROMAGAN TRANSFER AND RECEPTION OF RUSSAM RECEIPTION OF RUSSAM RECEIP					• Lagropartic Admin •			
Objects and Transfers								
Objects Object Transfers Routes	Obje	ts				Q Seard	1	
Sources	Showin	ng: All						
Archival Sources Books Newspapers and Periodicals/Reviews	C ritter Table Showing 10 v entries							/ ♥ entries
	-	Object name (Ricontrans) 🔶	Originator of Reference	Current \$	Archive 🖨	Creator 🖨	Card Status	ыф
Web Sources		Chalice in Margarites parish in Crete	/Organization/225, Parish of Margarites village, Rethymno	/Location/81, Margarites, Rethymno	<b>1</b>	katopi	unpublished	Object/10
Related Bibliography								
Passages and Comments		Icon of the Resurrection and Descent to the Hell with Twelve Great Feasts	/Organization/146, Byzantine museum Kastoria	/Location/79, Byzantine museum Kastoria	甗	maltezak	unpublished	Object/100
Source Passages Collection of Source Passages		Gospel of the Greek community of Kallipoli with Russian gilt silver cover	/Organization/2, Benaki Museum		18 2	spetridis	unpublished	Object/101
Researcher Comments Related Documentation	0	Icon of Sts. George and Demetrios on horseback	/Organization/148, Byzantine museum Kastoria	/Location/79, Byzantine museum Kastoria	of a	maltezak	unpublished	Object/102
Historical Figures Collections Events Locations Persons Organizations		Triptych depicting the Hospitality of Abraham, the Last Supper, the Marriage at Cana and the Twelve Great Feasts	/Organization/2, Benaki Museum		<b>686</b>	spetridis	unpublished	Object/103
	0	Icon with the Resurrection and Feasts	/Organization/2, Benaki Museum			spetridis	unpublished	Object/104