

U

Closed-class keywords
and
corpus-driven discourse analysis

B

Nicholas Groom
Centre for English Language Studies

Discourses

- **Sets of meanings and values which are associated with particular communities or institutions, and which are produced and reproduced through characteristic and often highly conventionalised linguistic choices (e.g. Carter 1995; Stubbs 1996, 1997, 2001; Partington 2003; Baker and McEnery 2005; Teubert 2005; Baker 2006; Mautner 2007).**

Corpora and discourses

- **Corpora allow for more robust statements to be made about discourses than can be made through the analysis of single texts**
- **Corpus analysis allows researcher to identify meanings and values in the language of many individuals in a discourse community**
 - Dominant and contested meanings
 - Synchronic and diachronic perspectives

Keywords, corpora and discourses

□ **Keywords analysis:**

- Delegates task of identifying objects of analysis to computer
- Not based on any particular theory of language
- Identifies items that would be difficult or impossible for human analyst to observe or predict
- Tends to generate very long lists of candidate items for analysis

Selecting keywords for discourse analysis

□ Strategy 1: Topslicing

- Words at top of list have highest sig. scores
- Ad hoc solution
- What lurks beneath?

Selecting keywords for discourse analysis

□ Strategy 2: Pruning

–3 kinds of keyword (Scott 1999; Scott & Tribble 2006):

□ proper nouns

□ other open-class ‘lexical’ words

□ closed-class ‘grammatical’ words

–Default strategy: discard or gloss over proper nouns and closed class words, and focus on remaining open-class items

Selecting keywords for discourse analysis

□ Strategy 2: Pruning

–Rationale:

- proper nouns only reveal *dramatis personae* and other obvious content features
- closed-class words only tell us about the style of a text or corpus of texts
- open-class words indicate the aboutness of a text or corpus, and are thus “generally those which are most interesting to analyse” (Baker 2006: 127).

Selecting keywords for discourse analysis

□ Strategy 2: Pruning

– Alternative strategy (Gledhill 2000; Groom 2007, forthcoming):

□ discard all open-class items and focus exclusively on remaining closed-class words

– Validity?

– Viability?

– Desirability?

Closed-class keywords as valid objects of analysis

□ **How meaningful are open-class words?**

Closed-class keywords as valid objects of analysis

□ How meaningful are open-class words?

– *gift* (n)

Closed-class keywords as valid objects of analysis

□ How meaningful are open-class words?

| | |
|---|--|
| You do have a <u>gift</u> for understatement, Jill. | |
| England gave Romania a <u>gift</u> of a goal. | |
| I wish I could [help you], but it's not in my <u>gift</u> . | |

Closed-class keywords as valid objects of analysis

□ How meaningful are open-class words?

| | |
|---|------------------------|
| You do have a <u>gift</u> for understatement, Jill. | talent |
| England gave Romania a <u>gift</u> of a goal. | unmissable opportunity |
| I wish I could [help you], but it's not in my <u>gift</u> . | power/authority |

Closed-class keywords as valid objects of analysis

□ How meaningful are open-class words?

| | |
|--|------------------------|
| HAVE + a + gift + for + N. | talent |
| a + gift + of + a + N | unmissable opportunity |
| N/PRON + BE + (not) + in + poss. det. + gift . | power/authority |

Closed-class keywords as valid objects of analysis

□ How meaningless are closed-class words?

| | |
|--|------------------------|
| <u>a</u> + gift + <u>for</u> + N. | talent |
| <u>a</u> + gift + <u>of</u> + <u>a</u> + N | unmissable opportunity |
| <u>in</u> + <u>poss. det.</u> + gift. | power/authority |

Closed-class keywords as valid objects of analysis

- **Closed-class words act as meaning classifiers (cf. pattern grammar, construction grammar)**

| | |
|--|------------------------|
| <u>a</u> + gift + <u>for</u> + N. | talent |
| <u>a</u> + gift + <u>of</u> + <u>a</u> + N | unmissable opportunity |
| <u>in</u> + <u>poss. det.</u> + gift. | power/authority |

Closed-class keywords as valid objects of analysis

- Meaning typically resides in *sequences* of words, not individual word forms
- Both open and closed-class words contribute to the meaning of an idiomatic sequence/pattern/construction/ ...
- ... so closed-class words are just as valid as open-class words are as starting points for semantically-oriented analysis
- Whither the style vs. aboutness distinction?

Closed-class keywords as viable objects of analysis

- **of**
- **very significant keyword in 3.2m-word corpus of history journal articles (HistArt) and 4m-word corpus of literary criticism journal articles (LitArt)**
- **'Style': nominalisation**

| Sequences | Examples | Sample 1 (%) | Sample 2 (%) | Sample 3 (%) |
|-------------------------------|---|--------------|--------------|--------------|
| <i>n of n</i> | it could not be truly democratic if it presupposed <u>the continuation of colonial control</u> . | 90 | 92 | 87 |
| <i>prep n</i> | they were eventually forced to seek help from local people, particularly from farmers <u>in the vicinity of the camps</u> . | 3 | 5 | 9 |
| <i>adj of n</i> | most of us are <u>aware of rules of evidence</u> | 1 | 2 | 2 |
| <i>v n out of n</i> | the laws acted to <u>take weapons out of men's hands</u> | 2 | 1 | 0 |
| <i>v of n</i> | <u>Bainville</u> ... <u>wrote of this recurrent dilemma of French foreign policy</u> | 1 | 0 | 2 |
| <i>fixed phrase</i> | they facilitated the entry of respectable women into what <u>one turn-of-the-century writer</u> termed the "Night Side of London" | 1 | 0 | 0 |
| <i>the adj-superl of pl-n</i> | the dubbing of a knight is <u>the most familiar of the new ceremonies</u> | 1 | 0 | 0 |
| <i>v n of n</i> | <u>Chulaki</u> <u>accused them both of cosmopolitanism</u> | 1 | 0 | 0 |

Closed-class keywords as viable objects of analysis

□ **of**

□ **very significant keyword in 3.2m-word corpus of history journal articles (HistArt) and 4m-word corpus of literary criticism journal articles (LitArt)**

□ **'Style': nominalisation**

□ **'Aboutness' ...?**

| Semantic sequence | Example | Sample 1 (%) | Sample 2 (%) | Sample 3 (%) | Average |
|--|---|--------------|--------------|--------------|---------|
| PROCESS + <i>of</i> + OBJECT | <u>control</u> of female sexuality | 11 | 17 | 15 | 14.333 |
| PROPERTY + <i>of</i> + PHENOMENON | the basic tenets of Marxist theory | 14 | 13 | 16 | 14.333 |
| CONCEPTUALISATION + <i>of</i> + PHENOMENON | a tidal wave of conservative <u>loyalism</u> | 7 | 15 | 15 | 12.333 |
| QUANTITY + <i>of</i> + PHENOMENON | Hundreds of thousands of shares | 7 | 8 | 7 | 7.333 |
| PROCESS + <i>of</i> + ACTOR | <u>the</u> death of Henry II | 6 | 7 | 6 | 6.333 |
| AUTHORITY + <i>of</i> + DOMAIN | <u>the</u> third Duke of Northumberland | 9 | 4 | 4 | 5.666 |
| PART + <i>of</i> + WHOLE | <u>the</u> east wall of that tiny church | 3 | 5 | 8 | 5.333 |
| QUALITY + <i>of</i> + PHENOMENON | <u>the</u> insignificance of the English | 3 | 7 | 5 | 5.000 |
| TEXT + <i>of</i> + CONTENT | Journal of Imperial and Commonwealth History | 9 | 2 | 2 | 4.333 |
| GROUP + <i>of</i> + MEMBERS | a transatlantic 'community of ideals, interests and purposes' | 4 | 2 | 3 | 3.000 |

Closed-class keywords as viable objects of analysis

□ **'Aboutness':**

□ **CONCEPTUALISATION + *of* + PHENOMENON:**

LitArt>HistArt

□ **Agentive processes (e.g. PROCESS + *of* + OBJECT):**

HistArt>LitArt

Closed-class keywords as preferred objects of analysis

- **Tractable lists of items for analysis**
 - **Gledhill (2000): 39 CCKWs**
 - **Groom (2007): 26 CCKWs**

Closed-class keywords as preferred objects of analysis

- **Tractable lists of items for analysis**
 - **Gledhill (2000): 39 CCKWs**
 - **Groom (2007): 26 CCKWs**
- **CCKWs tend to be scattered throughout KW lists (when ranked by keyness) - alternative to random sampling**

Closed-class keywords as preferred objects of analysis

□ Coverage (1)

- The 26 CCKWs obtained for HistArt alone constitute 20.28% of the whole corpus
- *“the majority of text is made of the occurrence of common words in common patterns”* (Sinclair 1991: 108)
- So it is arguably preferable to select the commonest of these common words for analysis.

Closed-class keywords as preferred objects of analysis

□ Coverage (2)

– Different types of phraseology

Closed-class keywords as preferred objects of analysis

□ Coverage (2)

- Different types of phraseology
- concordance sample of *of*:

Closed-class keywords as preferred objects of analysis

□ Coverage (2)

– Different types of phraseology

– concordance sample of *of*:

□ lexical bundles

1. debate began in London, **in advance of** the 1907 City Council el
2. men?"³⁷ His defensiveness **on behalf of** the scholastic prerogati
3. nds to loom ever larger **in the eyes of** electors, whether their
4. y, on June 8, 1794. ⁶⁶ **In the face of** such events, the madonna
5. foreign member states. **As a result of** these rulings, and relat
6. ularly from farmers **in the vicinity of** the camps.³¹ Although th

Closed-class keywords as preferred objects of analysis

□ Coverage (2)

– Different types of phraseology

– concordance sample of *of*:

□ grammar patterns

1. tion and most of us are **aware of** rules of evidence, which of ne
2. ed, we would **know** very little **of** the context of these oaths of
3. m of money. The vast **majority of** cases, however, went before th
4. icular, the frequent **outbreak of** disease, which posed an ostens
5. erer's own silence is **typical of** a larger phenomenon. Before th
6. brother, he **wrote** repeatedly **of** his bewilderment at the eighte

Closed-class keywords as preferred objects of analysis

□ Coverage (2)

– Different types of phraseology

– concordance sample of *of*:

□ semantic sequences

| CONCEPTUALISATION | <i>of</i> | PHENOMENON |
|-----------------------------|-----------|----------------------------|
| the complementary component | of | cultural stewardship |
| masculine images | of | cigar smoke and stiff wool |
| the institution | of | Mothers' Day |
| the problem | of | punctuation |
| a tidal wave | of | conservative loyalism |

Closed-class keywords as preferred objects of analysis

□ Coverage (2)

– Different types of phraseology

– concordance sample of *of*:

□ semantic sequences: *“recurring sequences of words and phrases that may be very diverse in form and which are therefore more usefully characterised as sequences of meaning elements rather than as formal sequences”* (Hunston 2008: 271)

Semantic sequences

- “*Hammerfest is a thirty-hour ride by bus from Oslo*”
- “*Ntobeye is a two-hour ride by four wheel drive vehicle from the vast refugee camp at Ngara*”

Semantic sequences

- “*Hammerfest is a NUMBER+TIME+JOURNEY by bus from Oslo*”
- “*Ntobeye is a NUMBER+TIME+JOURNEY by four wheel drive vehicle from the vast refugee camp at Ngara*”

Semantic sequences

- “*Hammerfest* is a NUMBER+TIME+JOURNEY by bus from Oslo”
- “*Ntobeye* is a NUMBER+TIME+JOURNEY by four wheel drive vehicle from the vast refugee camp at Ngara”

Semantic sequences

- “**DESTINATION** is a NUMBER+TIME+JOURNEY by bus from Oslo”
- “**DESTINATION** is a NUMBER+TIME+JOURNEY by four wheel drive vehicle from the vast refugee camp at Ngara”

Semantic sequences

- “ **DESTINATION** is a NUMBER+TIME+JOURNEY by *bus* from Oslo”
- “ **DESTINATION** is a NUMBER+TIME+JOURNEY by *four wheel drive vehicle* from the vast refugee camp at Ngara”

Semantic sequences

- “ **DESTINATION** *is a* **NUMBER TIME JOURNEY** *by* **MODE OF TRANSPORT** *from Oslo*”
- “ **DESTINATION** *is a* **NUMBER TIME JOURNEY** *by* **MODE OF TRANSPORT** *from the vast refugee camp at Ngara*”

Semantic sequences

- “ **DESTINATION** *is a* **NUMBER TIME JOURNEY** *by* **MODE OF TRANSPORT** *from* *Oslo*”
- “ **DESTINATION** *is a* **NUMBER TIME JOURNEY** *by* **MODE OF TRANSPORT** *from* *the vast refugee camp at Ngara*”

Semantic sequences

- “DESTINATION *is a* NUMBER TIME JOURNEY *by* MODE OF TRANSPORT *from* POINT OF DEPARTURE”
- “DESTINATION *is a* NUMBER TIME JOURNEY *by* MODE OF TRANSPORT *from* POINT OF DEPARTURE”

Semantic sequences

- *“Hammerfest is a thirty-hour ride by bus from Oslo”*
- *“Ntobeye is a two-hour ride by four wheel drive vehicle from the vast refugee camp at Ngara”*

Semantic sequences

- “**DESTINATION** *is a* **NUMBER TIME JOURNEY** *by* **MODE OF TRANSPORT** *from* **POINT OF DEPARTURE**”
- “**DESTINATION** *is a* **NUMBER TIME JOURNEY** *by* **MODE OF TRANSPORT** *from* **POINT OF DEPARTURE**”

Semantic sequences

- “DESTINATION *is a* NUMBER TIME JOURNEY *by* MODE OF TRANSPORT *from* POINT OF DEPARTURE”
- “DESTINATION *is a* NUMBER TIME JOURNEY *by* MODE OF TRANSPORT *from* POINT OF DEPARTURE”

Conclusion

- **Closed-class keywords are**
 - **valid**
 - **viable**
 - **perhaps even preferable**
- for corpus-driven analyses of specialized discourses**

Coda

□ Student essay (2009):

– *“The grammatical words ‘of’ and ‘and’ cannot tell us anything at all.”*