

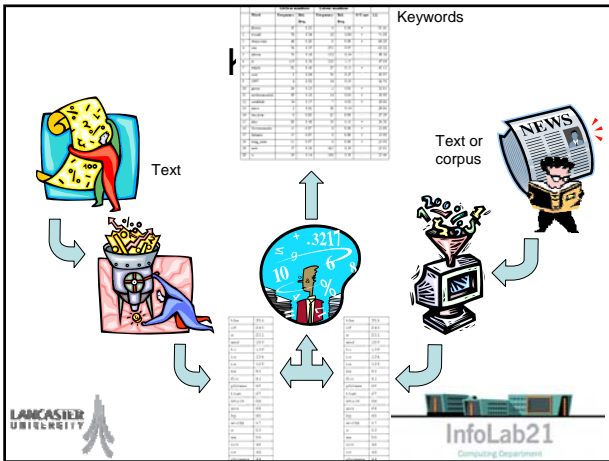
The key domain method for the study of language varieties

Paul Rayson and Nicholas Smith
UCREL, Computing Department,
Lancaster University, UK.



Outline

- New kind of method and tool (Matrix) for the statistical analysis of corpora
- Matrix integrates part-of-speech tagging and semantic field tagging in a profiling tool
- Extends the keywords procedure to identify key grammatical categories and key concepts
- Case studies:
 - Key concepts in 20th C. romantic fiction
 - 'Obligation & necessity' in 20th C. corpora
- Software demo



Log likelihood statistic

	Corpus one	Corpus two	Total
Frequency of word	a	b	a+b
Frequency of word not occurring	c-a	d-b	c+d-a-b
TOTAL	c	d	c+d

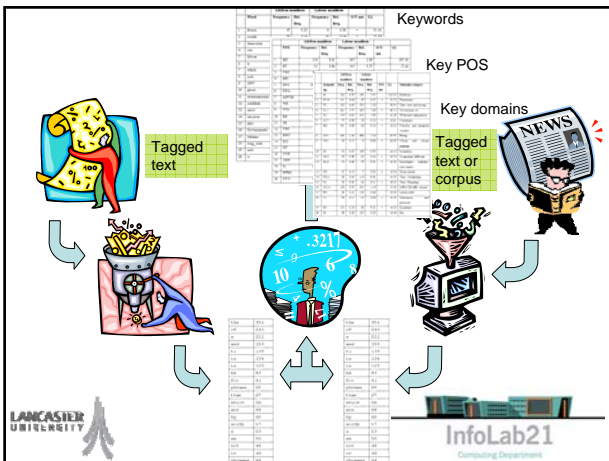
$$E_i = \frac{N_i \sum O_i}{\sum N_i}$$

$$-2 \ln \lambda = 2 \sum O_i \ln \left(\frac{O_i}{E_i} \right)$$

$$E1 = c \times (a+b) / (c+d)$$

$$E2 = d \times (a+b) / (c+d)$$

$$LL = 2 \times ((a \times \ln(a/E1)) + (b \times \ln(b/E2)))$$



Matrix method

- Integrates POS tagging and semantic field annotation into a profiling tool
- Extends keywords procedure to identify key grammatical categories and key concepts



Previous studies using Matrix

- Social differentiation in the use of English vocabulary (Rayson, Leech & Hodges, 1997)
- Profiling of learner English (Granger & Rayson, 1998)
- Semantic analysis of technical documents from the software engineering domain (Sawyer, Rayson & Garside, 2005)



The data for today's case studies

- Standard written British English
- Sampling dates 1931 – 1961 – 1991
- Each corpus contains 1M words
- 15 genres of published informative and imaginative prose (e.g. press reportage, academic writing, romantic fiction, science fiction)
- Corpus size and sampling frame modelled on the Brown Corpus (of 1960s American English)



Comparable corpora across the C20th

	1901 (1898-1904)	1931 (1928-34)	1961	1991/ 1992
BrE	Lanc-1901	B-LOB (Lanc-1931)	LOB	F-LOB
AmE	?	Pre-Brown31	Brown	Frown



Example 1: A bottom-up view of Romantic fiction across the 20th C.

General questions:

- What 'key concepts' help to distinguish romantic fiction from fiction in general?
- What changes among 'key concepts' have changed over the course of the 20th century?



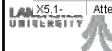
Romantic fiction vs. Fiction in general

- Synchronic comparison between the F-LOB corpus and the BNC-Sampler Imaginative subcorpus
- Intuitive expectations about key concepts: LOVE, RELATIONSHIPS, ...



Items used significantly more frequently in Romantic fiction than in general fiction (FLOB vs. BNC-Sampler Imaginative subcorpus)

Tag	Gloss	FLOB romantic		Sampler Imag		Over/under use	LL
		freq	%	freq	%		
A5.3+	Evaluation: Accuracy	125	0.2	150	0.1	+	69.0
T1.1	Time: General	27	0.0	4	0.0	+	58.3
L1+	Life and living things	47	0.1	38	0.0	+	42.1
X2.6-	Expect	40	0.1	28	0.0	+	40.9
B1	Anatomy / physiology	1104	1.7	3057	1.4	+	33.6
X2.6+	Expect	74	0.1	101	0.1	+	33.1
X2.4	Investigate, examine, test, search	128	0.2	233	0.1	+	30.1
S1.1.1	Social actions, states and processes	127	0.2	237	0.1	+	27.8
T1.3-	Time: Period	58	0.1	77	0.0	+	27.2
O4.2+	Judgement of appearance (pretty etc.)	194	0.3	448	0.2	+	19.3
N3.7+	Measurement: Length & height	86	0.1	164	0.1	+	17.7
O4.1	General appearance and physical properties	142	0.2	322	0.1	+	15.4
A12++	Easy/difficult	11	0.0	7	0.0	+	12.2
X5.1-	Attention	25	0.0	34	0.0	+	11.2



Key tag 'EXPECT' in Romantic fiction (1990s BrE)

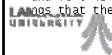
Her attitude took Frances by surprise around him, catching him by surprise most of Father's unexpected surprises. <S There was something unpredictable er nipple as she gasped with surprise y, when their future seemed hopeless rs. </p><p><S <quote> How e you ever known anything so </p><p><S <quote> I'm not discoveries. <S It would be Abby turned to look at him, ne. </p><p><S To Jenny's s To Jenny's surprise and f her could n't help feeling nger. <S She was constantly s that, as Thomas aged with s the driver turned in some o the driver, he spoke with e to you often. <S But I wo </p><p><S To her complete surprise, for it was in such marked c - which might not have matter, in which the gift was of su in her that was only explaine , sucking so it puckered poin , Judy had an idea. </p><p> quiet it is, </quote> she sa quiet? </quote> </p><p><S , </quote> said Tom. </p><p> if there were. <S I was hopi at her own sense of relief at and - astonishingly enough - enough - relief, he appeared that he had taken her tentati by the sensitive perception o rapidly, Sophie became even , then laughed to see a maid authority. </p><p><S <quot waiting for a letter from you , she received long letters o



InfoLab21
Computing Department

Key tag 'ATTENTION' in Romantic fiction (1990s BrE)

ad had their orders and they ignored <p><S His younger brother ignored low Jocelyn's behaviour to upset use and because he wanted to distract <quote> Did you feel she 's a castle walls. </p><p><S Ignoring proving quite devastatingly distracting . </p><p><S <quote> Do n't upset in a dream art I pushed open the door, ignoring put off asking Rob what was troubling Bob's voice broke into her reverse with Danny she would n't sit aimlessly in a dream am Foley was asleep and deep diversion take your mind off <quote> But that did n't take your mind off ignoring they had finished and then, upset and it. <S I do n't want to be ignored <quote> she 'll soon



InfoLab21
Computing Department

Diachronic development of key concepts in Romantic fiction (1): Items used significantly more frequently in 1991 than in 1961

Tag	Gloss	FLOB 1961		Sampler 1991		Over/under use	LL
		freq	%	freq	%		
B1	Anatomy and physiology	1104	1.7	756	1.3	*	24.7
B3	Medicines and medical treatment	131	0.2	53	0.1	*	23.8
K4	Drama, the theatre and showbusiness	88	0.1	25	0.0	+	28.6
X2.2*	Sensory: Sound	16	0.0	3	0.0	+	8.0
S1.2.5*	Toughness, strong/weak	25	0.0	8	0.0	+	6.9



InfoLab21
Computing Department

Key concept 'ANATOMY/PHYSIOLOGY' in 1991 (FLOB)

shattered. <S She could n't sleep eyes with their big o s appeared around her lovely quote> </p><p><S Sophie's eyes filled with tears <quote> Oh, Thomas, yo was behaving itself. <S He p with a lean brown hand, <S <S <quote> Of course I want only by willpower and a tight was dazzling and his curly ye showed off his tan. </p><p> </p><p><S <quote> Well, her cheeks slightly spotty slightly spotty. <S There ha appointments, and the glossy for a few days on a low-starc <S He's kept the beard, t , though. <S He says it 's g </p><p><S THE CHAPEL outing on the Perkins girl, Ann. <



InfoLab21
Computing Department

Diachronic development of key concepts in Romantic fiction (2): Items used significantly more frequently in 1961 (LOB) than in 1931 (B-LOB)

Tag	Gloss	FLOB 1961		B-LOB 1931		Over/under use	LL
		freq	%	freq	%		
T1.1.3	Time: General: Future	463	0.8	346	0.6	*	24.9
G2.1-	Crime, law and order.	54	0.1	25	0.0	*	12.7
H2	Parts of buildings	245	0.4	186	0.3	*	12.1
G2.1	Crime, law and order.	50	0.1	23	0.0	+	11.9
X2.6*	Expect	86	0.2	51	0.1	+	11.3
Z1	Personal names	1506	2.7	1117	1.9	+	84.1
Z3	Other proper names	117	0.2	71	0.1	+	14.4
M1	Moving, coming and going	909	1.61	841	1.4	+	8.42
W2	Light	73	0.13	46	0.08	+	7.94
A7*	Definite (* modals)	705	1.25	649	1.08	+	6.99
A14	Exclusivizers/particularizer s	254	0.45	212	0.35	+	6.79
Q2.2	Speech acts	665	1.18	611	1.02	+	6.72
Z8	Pronouns etc.	7436	13.16	7577	12.62	+	6.68



InfoLab21
Computing Department

Key concept 'TIME: GENERAL: FUTURE' in 1961 (LOB) vs. 1931 (B-LOB)

<quote> ooh, it 's late - I 'll going home for my dinner or I 'll Tommy's direction. <quote> I 'll 'cos Grandma 's ill but she 's going to soon the Town Hall clock when it 's going to in across the road. <quote> I 'll she did n't mind, in fact the sooner the better and more luck next time <quote> oh that - I expect you 'll te> one of your many young men will on. <quote> come along and we 'll from Glasgow to talk over her future pson sounded anxious. <quote> will is n't it kind of her? I - I 'll fit and - <quote> <quote> you will ave lunch with her at her flat <quote> she explained, <quote> she 'll soon have to be going home for my d be late back for school <quot see you across the road get a job soon </quote>. he p </quote>. he paused to consid strike </quote>. he hurried se take you home. </quote> he co the better and more luck next was her motto. she had been t be thinking of having one like be sweeping you off your feet have a cup of tea before I tak with Mr Robertson and the loca n't you be lonely? </quote> < never forget all she did for m need a full purse! </quote> M . then she followed the maid u soon make you feel at home. < make you feel at home. </quot



InfoLab21
Computing Department

Diachronic development of key concepts in Romantic fiction (3): Items used significantly more frequently in 1931 (B-LOB) than in 1961(LOB)

Tag	Gloss	B-LOB 1931		LOB 1961		Over/under use	LL
		freq	%	freq	%		
A13.5	Degree: Compromisers	104	0.2	49	0.1	+	17.0
A13.2	Degree: Maximizers	124	0.2	64	0.1	+	16.0
E4.1+++	Happy/sad: Happy	17	0.0	2	0.0	+	12.7
N3.2-	Measurement: Size	156	0.3	95	0.2	+	11.5
S7.1-	Power, organizing	46	0.1	18	0.0	+	11.0
G1.2	Politics	20	0.0	2	0.0	+	16.0
A1	GENERAL AND ABSTRACT TERMS	8	0.0	0	0.0	+	10.6
N1	Numbers	341	0.6	245	0.4	+	10.5
A5.1+	Evaluation:- Good/bad	169	0.3	108	0.2	+	10.1
A13.3	Degree: Boosters	417	0.7	313	0.6	+	9.2
A13.7	Degree: Minimizers	65	0.1	33	0.1	+	8.8
T3+	Time: Old, new and young, age	176	0.3	117	0.2	+	8.6



Key concept 'DEGREE: COMPROMISERS' in 1931 (B-LOB) vs. 1961 (LOB)

crouch on the floor and stay a moment she would come back, a mile away." "That is closely viewed, he might be ake another in order to make them, seemed to have drunk le there, who were not only alf hoped that something not turned out afterwards to be ff' and I certainly found it. Perhaps I even swaggered, back into her eyes and face ge into the hotel. She felt at this hotel before. It's uppose you're pleased. And ng of their own later on." at would n't die, he knew, e ways. It came to him with

quite quiet quite quiet quite quiet rather rather rather rather rather rather rather rather rather

still, and if the worst sho confident to say that we sto right," says she: "but y a plain young man. His eyes sure that I had not been dec enough, and they were talki quiet and subdued but looked proper might happen at any m expensive. She sipped her g exhilarating myself. The di that I scarcely remember in steadily and with an express queer as Harvey signed the v jolly, is n't it?" But Ja superior about it, into the nice," he said. "In fact unhappily, that he was afra a shock, once he'd kissed



Example 2: Key concept OBLIGATION/NECESSITY in 20th C. BrE

Background:

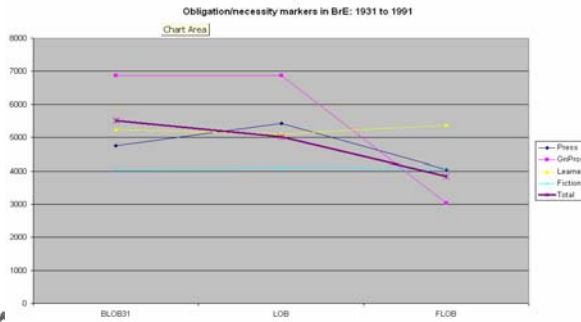
- Recent observations of significant shifts having occurred among expressions of obligation/necessity in the period 1961-1991 (Leech 2003, Smith 2003). E.g.
 - a decline of the central modals MUST and NEED
 - a spread of the semi-modals HAVE TO, NEED TO

Questions

- Are these changes recent?
- How do these changes compare to the development of the semantic field of OBLIGATION/NECESSITY as a whole?



Example 2: Key concept OBLIGATION/NECESSITY in 20th C. BrE

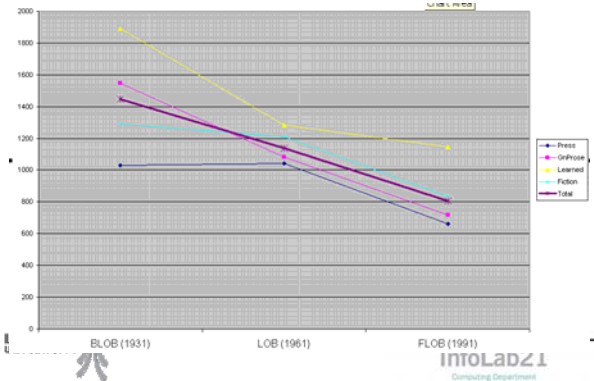


'OBLIGATION/NECESSITY' in 1931 Fiction

ugh the sight reassured me, I must confess that I dreaded taking p ut whether he was or was not I should not know till one of two things e her, I began to fear that I must have passed her by: then the d aid I, "I shall hear all you have to say: and you are to sit well b ordward whatever befalls. If I should stop or be stopped, you are to quite still, and if the worst should happen, you must happen, you must swear that yo the worst should happen, you must swear that you do not know me, have known that Lelia, to whom om roads meant nothing, I should see about her if she was to tel have to have to tel I not take you further. As it is home till long after dark. It must be six miles to Merring from wh o see her. You think you will present me and that I am not fi tter," said I. "And now you must go to your home. Soon after da you do not come back?" "Why go about your business," said you id you good-bye to-day." "Why you?" said I. "If I do not c towered her eyes. "I do not see her," she said. The l n't." "Jenny said, with decs have had a startling beauty in o go home." "I'll go. You need a close inspection. Kings did tention, but Jenny's charms had given that thus her father must once have looked, and her own



MUST in BrE: 1931 to 1991



Conclusion (1)

- Matrix extends keywords approach to key grammatical classes and key concepts using Log Likelihood to compare frequency profiles
- Key grammatical categories and semantic classes are used to group together lower frequency words and those words which would, by themselves, not be identified as key, and would otherwise be overlooked
- Comparison at the POS and semantic levels reduces the number of key categories that the researcher needs to examine



Conclusion (2): Planned future developments

- Semantic tagging:
 - More accurate semantic disambiguation
 - A more refined semantic analysis for certain categories, e.g. deontic vs. epistemic use of MUST, HAVE TO etc.
- Semantic profiling
 - Incorporation of a Dispersion measure to filter 'key concepts' tables
 - Option to filter text according to markup types, e.g. quoted speech



Further info

- <http://ucrel.lancs.ac.uk/20thCenturyEnglish/>
- Contacts: Paul Rayson (paul@comp.lancs.ac.uk), Nick Smith (nick@comp.lancs.ac.uk)
- Sponsored by The Leverhulme Trust (Grant number F/00 185/J), this project runs from August 2005 - July 2007.



References

- Rayson, P. (2003). Matrix: A statistical method and software tool for linguistic analysis through corpus comparison. *Ph.D. thesis*, Lancaster University.
- Rayson, P., Leech, G., and Hodges, M. (1997). Social differentiation in the use of English vocabulary: some analyses of the conversational component of the British National Corpus. *International Journal of Corpus Linguistics*, Volume 2, number 1, pp 133 - 152. John Benjamins, Amsterdam/Philadelphia.
- Granger, S., and Rayson, P. (1998). Automatic profiling of learner texts. In S. Granger (ed.) *Learner English on Computer*. Longman, London and New York, pp. 119-131.
- Sawyer, P., Rayson, P. and Cosh, K. (2005) Shallow Knowledge as an Aid to Deep Understanding in Early Phase Requirements Engineering. *IEEE Transactions on Software Engineering*, Volume 31, number 11, November, 2005, pp. 969 - 981.

