

The advantages of using relational databases for large historical corpora

Workshop on Historical Text Mining / U Lancaster / July 2006

1. VIEW: Variation in English Words and Phrases (<http://view.byu.edu>)

1.1 Original format:

<w AT0>the <w NN1>reluctance <w PRF>of <w AJ0>central <w NN1>government
<w TO0>to <w VVI>take <w AVP>on <w DT0>any <w AJ0>new <w NN1>expenditure

1.2 Convert to "vertical format" and add delimiters

AT0^the
NN1^reluctance
PRF^of
AJ0^central
NN1^government
TO0^to

Table 1. Import into MS SQL Server ; run SQL query to add adjacent words

ID	word1	pos1	word2	pos2	word3	pos3	word4	pos4	word5	pos5	word6	pos6	word7	pos7
1463998	the	AT0	reluctance	NN1	of	PRF	central	AJ0	government	NN1	to	TO0	take	VVI
1463999	reluctance	NN1	of	PRF	central	AJ0	government	NN1	to	TO0	take	VVI	on	AVP
1464000	of	PRF	central	AJ0	government	NN1	to	TO0	take	VVI	on	AVP	any	DT0
1464001	central	AJ0	government	NN1	to	TO0	take	VVI	on	AVP	any	DT0	new	AJ0
1464002	government	NN1	to	TO0	take	VVI	on	AVP	any	DT0	new	AJ0	expenditure	NN1
1464003	to	TO0	take	VVI	on	AVP	any	DT0	new	AJ0	expenditure	NN1	makes	VVZ
1464004	take	VVI	on	AVP	any	DT0	new	AJ0	expenditure	NN1	makes	VVZ	it	PNP
1464005	on	AVP	any	DT0	new	AJ0	expenditure	NN1	makes	VVZ	it	PNP	necessary	AJ0
1464006	any	DT0	new	AJ0	expenditure	NN1	makes	VVZ	it	PNP	necessary	AJ0	to	TO0
1464007	new	AJ0	expenditure	NN1	makes	VVZ	it	PNP	necessary	AJ0	to	TO0	seek	VVI
1464008	expenditure	NN1	makes	VVZ	it	PNP	necessary	AJ0	to	TO0	seek	VVI	funds	NN2

Table 1. Basic query: break * [n*]

DISTRIB	WORD/PHRASE	TOKENS REG1	PER MIL IN REG1 [100,000,000 WORDS]
1	BREAK THE LAW	119	1.19
2	BREAK THE NEWS	58	0.58
3	BREAK THE DEADLOCK	41	0.41
4	BREAK THE RULES	31	0.31
5	BREAK THE ICE	29	0.29
6	BREAK THE SILENCE	28	0.28
7	BREAK WITH TRADITION	28	0.28
8	BREAK THE HABIT	24	0.24

Table 2. By register (e.g. lexical verbs in [law texts] (W_ac_polit_law_edu))

DISTRIB	WORD/PHRASE	TOKENS REG1	PER MIL IN REG1 [4,640,346 WORDS]
1	MAKE	3133	675.17
2	TAKE	2276	490.48
3	GIVE	1792	386.18
4	PROVIDE	1347	290.28
5	APPLY	1278	275.41
6	SEE	1212	261.19
7	USE	1137	245.02
8	CONSIDER	1041	224.34

Table 3. Specific to register (e.g. compare to other ACAD) (x ≥ 2 in both)

DISTRIB	WORD/PHRASE	TOKENS REG1	TOKENS REG2	PER MIL IN REG1 [4,640,346 WORDS]	PER MIL IN REG2 [10,789,236 WORDS]	REG 1-2 RATIO
1	SUE	331	10	71.33	0.93	76.96
2	CERTIFY	27	2	5.82	0.19	31.39
3	ADJOURN	26	2	5.60	0.19	30.23
4	NOTIFY	38	3	8.19	0.28	29.45
5	WAIVE	38	3	8.19	0.28	29.45
6	DISCLOSE	190	16	40.95	1.48	27.61
7	OVERRULE	23	2	4.96	0.19	26.74
8	PROHIBIT	53	5	11.42	0.46	24.65
9	PLEAD	62	6	13.36	0.56	24.03
10	RESCIND	29	3	6.25	0.28	22.48

Table 4. See relative frequency in all 70 registers (*plead*)

#	REGISTER NAME	# PER MILLION	# TOKENS	# WORDS
1	S courtroom	102.0	13	127474
2	W fict drama	43.7	2	45757
3	W letters_prof	15.1	1	66031
4	W ac polit law edu	14.2	66	4640346
5	W fict poetry	13.5	3	222451
6	W hansard	11.2	13	1156171
7	W news script	10.8	14	1292156
8	S parliament	10.4	1	96239
9	W_email	9.4	2	213045
10	W newsp other arts	8.4	2	239258
11	S_consult	7.2	1	138011
12	W newsp brdsht nat misc	5.8	6	1032943
13	W newsp brdsht nat arts	5.7	2	351811
14	W newsp tabloid	5.5	4	728413
15	W religion	5.3	6	1121632

Table 5. Lexical bundles (ACAD vs FICT)

DISTRIB	WORD/PHRASE	TOKENS REG1	TOKENS REG2	PER MIL IN REG1 [15,429,582 WORDS]	PER MIL IN REG2 [16,194,885 WORDS]	REG 1-2 RATIO
1	THE USE OF	3179	170	206.03	10.50	19.63
2	AS A RESULT	2063	143	133.70	8.83	15.14
3	IN WHICH THE	2589	181	167.79	11.18	15.01
4	THE FORM OF	1487	116	96.37	7.16	13.45
5	THE NUMBER OF	2701	213	175.05	13.15	13.31
6	THE UNITED STATES	1768	140	114.59	8.64	13.25
7	THE EFFECT OF	1662	133	107.72	8.21	13.12
8	THAT THERE IS	1635	137	105.97	8.46	12.53
9	CAN NOT BE	2858	242	185.23	14.94	12.40

Table 6. Collocates: [*kitchen*] within 10 words of a noun (by frequency)

	WORD	# TIMES NEARBY	TOTAL IN CORPUS	% NEARBY
1	DOOR	518	23207	2.2%
2	ROOM	460	28443	1.6%
3	TABLE	425	18804	2.3%
4	HOUSE	261	47087	0.6%
5	BATHROOM	205	2323	8.8%

Table 7. [*kitchen*] within 10 words of a noun (by percentage / modified Z-score)

	WORD	# TIMES NEARBY	TOTAL IN CORPUS	% NEARBY
1	DRAINER	12	32	37.5%
2	SINK	92	285	32.3%
3	SCULLERY	32	176	18.2%
4	WORKTOP	13	75	17.3%
5	UTENSILS	22	127	17.3%
6	LIVING-ROOM	11	69	15.9%
7	PANTRY	18	155	11.6%

Table 8. CHARTS: frequency of [nn*] [nn*] – six “macro” registers

REGISTER	<u>SPOKEN</u> 4887	<u>FICTION</u> 6076	<u>NEWS</u> 11362	<u>ACADEMIC</u> 16485	<u>NON-ACAD</u> 18960	<u>MISC</u> 34146
	472.9	375.2	1,068.1	1,068.4	1,139.8	1,202.6

Table 9. CHARTS: Frequency of [nn*] [nn*] – sub-registers of SPOKEN

	DIS_	DOC_	NWS	CLS_	CST_	CNV_	CRT_	DEM_	INT_	ORH_	COM_	HUM_	NAT_	POL_	SOC_	MTG_	PRL_	DEB_	SRM_	SP+_	SP_-	SPO_	TUT_	UNC
	515	650	689	291	312	305	604	283	509	450	1,390	118	1,235	413	375	744	873	959	109	934	660	270	615	486

Table 10. Grouping by synonyms: {sheer/utter/absolute} [nn*]

	%	+	-	SHEER		%	+	-	UTTER		%	+	-	ABSOLUTE
1	1.00	60	--	<u>weight</u>	1	1.00	19	--	<u>confusion</u>	1	1.00	98	--	<u>majority</u>
2	1.00	31	--	<u>force</u>	2	1.00	5	--	<u>condemnation</u>	2	1.00	84	--	<u>terms</u>
3	1.00	26	--	<u>luck</u>	3	1.00	5	--	<u>devastation</u>	3	1.00	54	--	<u>zero</u>
4	1.00	23	--	<u>quantity</u>	4	1.00	5	--	<u>disregard</u>	4	1.00	53	--	<u>minimum</u>
5	1.00	13	--	<u>cliff</u>	5	1.00	5	--	<u>helplessness</u>	5	1.00	51	--	<u>value</u>
6	1.00	12	--	<u>cliffs</u>	6	1.00	4	--	<u>loneliness</u>	6	1.00	39	--	<u>egalitarianism</u>
7	1.00	11	--	<u>coincidence</u>	7	1.00	3	--	<u>dejection</u>	7	1.00	33	--	<u>right</u>
8	1.00	10	--	<u>enjoyment</u>	8	1.00	3	--	<u>ruthlessness</u>	8	1.00	27	--	<u>price</u>

Table 11. Comparing collocates: man/woman + ADJ

	MAN +	# MAN	# WOMAN	% MAN		WOMAN +	# WOMAN	# MAN	% WOMAN
2	<u>UNITED</u>	48	0	100.0%	1	<u>CLEANING</u>	28	0	100.0%
3	<u>MACHO</u>	36	0	100.0%	2	<u>DUMPY</u>	23	0	100.0%
4	<u>BURLY</u>	30	0	100.0%	4	<u>LIBERATED</u>	18	0	100.0%
5	<u>SELF-MADE</u>	29	0	100.0%	5	<u>MOTHERLY</u>	17	0	100.0%
6	<u>ARMED</u>	26	0	100.0%	11	<u>BUXOM</u>	11	0	100.0%
10	<u>CIVILISED</u>	19	0	100.0%	12	<u>RAPED</u>	10	0	100.0%
11	<u>MUSCULAR</u>	19	0	100.0%	13	<u>VAGINAL</u>	8	0	100.0%
12	<u>SCIENTIFIC</u>	17	0	100.0%	20	<u>ABUSED</u>	6	0	100.0%
13	<u>STEADY</u>	16	0	100.0%	23	<u>BOSSY</u>	6	0	100.0%
16	<u>TALENTED</u>	14	0	100.0%	24	<u>CLINGING</u>	6	0	100.0%
21	<u>CIVILIZED</u>	12	0	100.0%	25	<u>COMPLAINING</u>	6	0	100.0%

Table 12. Comparing collocates: chair in FICTION and ACADEMIC

	WORD	# REG 1	# REG 2	% REG 1		WORD	# REG 2	# REG 1	% REG 2
1	<u>HANDS</u>	84	0	100.0%	1	<u>COMMITTEE</u>	18	0	100.0%
2	<u>LEGS</u>	62	0	100.0%	2	<u>FIELD</u>	17	0	100.0%
3	<u>FIRE</u>	59	0	100.0%	3	<u>MEMBERSHIP</u>	7	0	100.0%
4	<u>FEET</u>	58	0	100.0%	4	<u>PHILOSOPHY</u>	7	0	100.0%
5	<u>KITCHEN</u>	55	0	100.0%	5	<u>CIRCLE</u>	4	0	100.0%
6	<u>FACE</u>	50	0	100.0%	6	<u>CO-ORDINATOR</u>	4	0	100.0%
7	<u>FLOOR</u>	47	0	100.0%	7	<u>REFERENCE</u>	4	0	100.0%
8	<u>WOMAN</u>	31	0	100.0%	8	<u>ROLE</u>	4	0	100.0%
9	<u>END</u>	29	0	100.0%	9	<u>MATHEMATICS</u>	3	0	100.0%
10	<u>LEATHER</u>	29	0	100.0%	10	<u>FRENCH</u>	2	0	100.0%
11	<u>COFFEE</u>	28	0	100.0%	11	<u>ADMISSIONS</u>	2	0	100.0%
12	<u>JACKET</u>	28	0	100.0%	12	<u>INJURIES</u>	2	0	100.0%

Table 13. WordNet queries

SYNONYM: [=x] [=small] [=scream].[v*]	MORE GENERAL: [>x] [>shriek] [>wallop].[v*]	MORE SPECIFIC: [<x] [<woman] [<hit].[v*]	PART OF: [@x] [@house] [@body]	CONTAINS PART: [&x] [&leg] [&shelf]
--	--	---	---	--

Table 14: [=beat].[v*] the [nn*]

	%	+	-	BEAT		%	+	-	CRUSH		%	+	-	VANQUISH
1	1.00	16	--	system	1	1.00	4	--	rebellion	1	1.00	1	--	alien
2	1.00	14	--	eggs	2	1.00	2	--	biscuits	2	1.00	1	--	romans
3	1.00	14	--	world	3	1.00	2	--	coup	3	0.50	1	1	past
4	1.00	11	--	recession	4	1.00	2	--	imagination					
5	1.00	10	--	egg	5	1.00	1	--	amaretti					
6	1.00	7	--	hell	6	1.00	1	--	bottle					
7	1.00	6	--	clock	7	1.00	1	--	cane					
8	1.00	6	--	cream	8	1.00	1	--	career					
9	1.00	6	--	shit	9	1.00	1	--	cockatrice					
10	1.00	6	--	traffic	10	1.00	1	--	corn					

Table 15: More specific words for walk as verb: [<walk].[v*]

18	STROLL	485	4.85
19	LIMP	466	4.66
20	TREAD	445	4.45
33	STALK	224	2.24
35	LURCH	201	2.01
36	STUMBLE	193	1.93
39	TIPTOE	144	1.44
41	STAGGER	128	1.28
47	STRUT	95	0.95
49	SWAGGER	77	0.77
54	PROWL	60	0.60
55	TRUDGE	60	0.60
62	SAUNTER	34	0.34
67	TRAIPISE	18	0.18
68	PRANCE	18	0.18