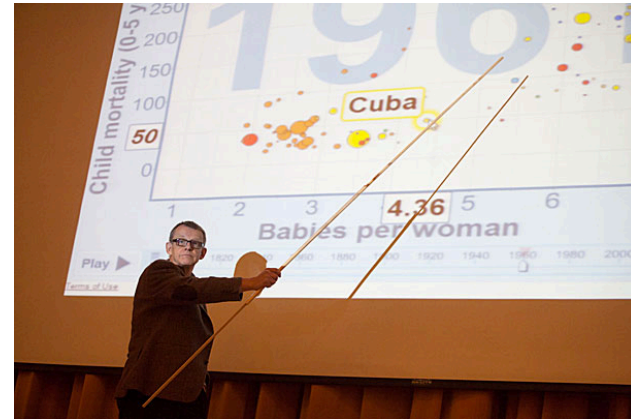


Analyzing linguistic change with motion charts

Martin Hilpert

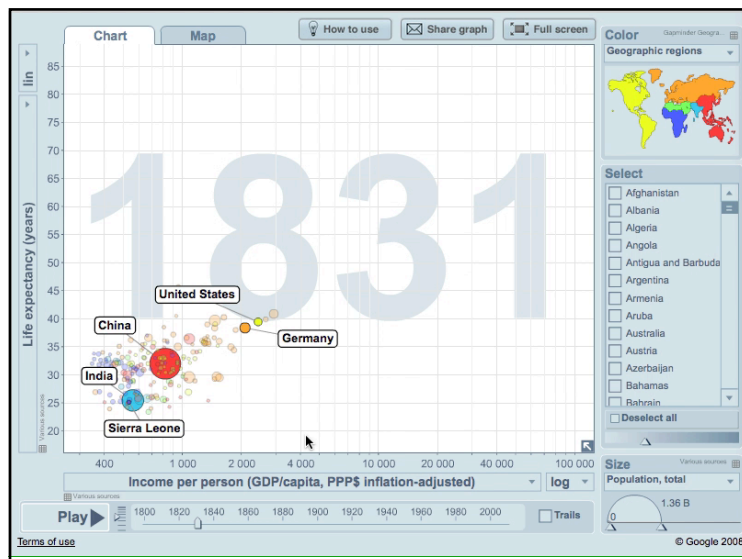


1



Hans Rosling

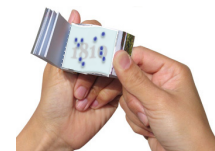
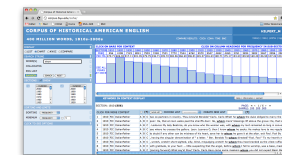
2



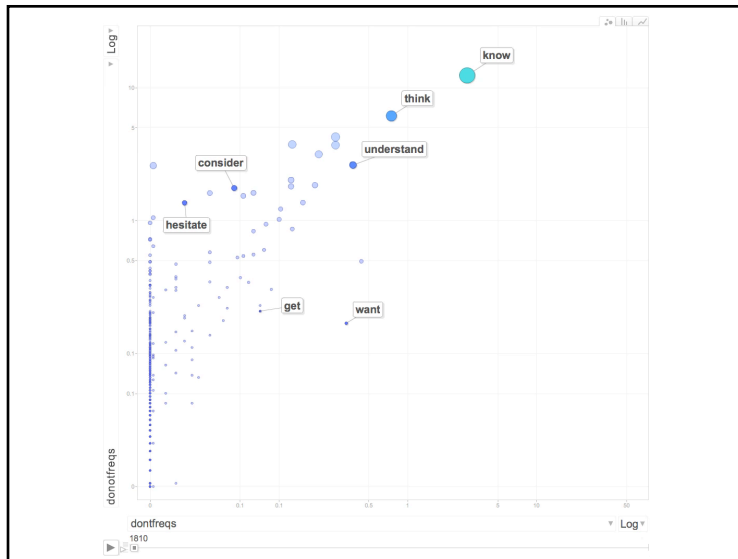
3

What about linguistic data?

- Use a corpus that represents identical kinds of text across multiple periods of time
- Select a phenomenon and create a visualization for each corpus period.
- View the visualizations in sequence.



4



5

Plan for today

- Two brief case studies
 1. noun-participle compounds (*whisky-soaked*)
 2. many a NOUN (*many a surprise*)
- Making your own motion charts

6

1. Noun-participle compounds

7

Huddleston and Pullum (2002)

- *drug-related, home-made, safety-tested, taxpayer-funded*
- “These compounds generally correspond to syntactic passives with a PP: *related to drugs, made at home, tested for safety, funded by taxpayers, etc.*” (2002: 1659)

8

Bauer et al. (2013)

- The first element “cannot receive an object interpretation”
 - doctor-recommended (SUBJECT)
 - arsenic-exposed (PREPOSITIONAL OBJECT)
 - *lunch-eaten (OBJECT)



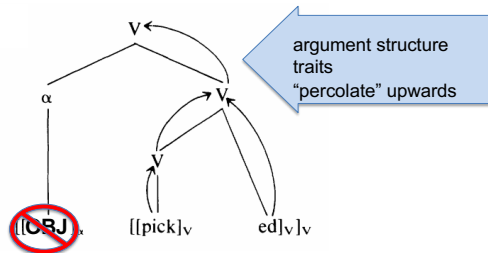
9

Lieber (1983)

- The past participle changes the argument structure of a transitive verb.
- The past participle ‘expels’ the direct object:
 - passive:
 - The object of a transitive verb cannot be part of the verb phrase.
 - *The strawberries [were picked by hand]_{VP}*
 - noun-participle compounding:
 - The object of a transitive verb cannot be part of the compound.
 - the [hand-picked]_A strawberries

10

Lieber (1983)

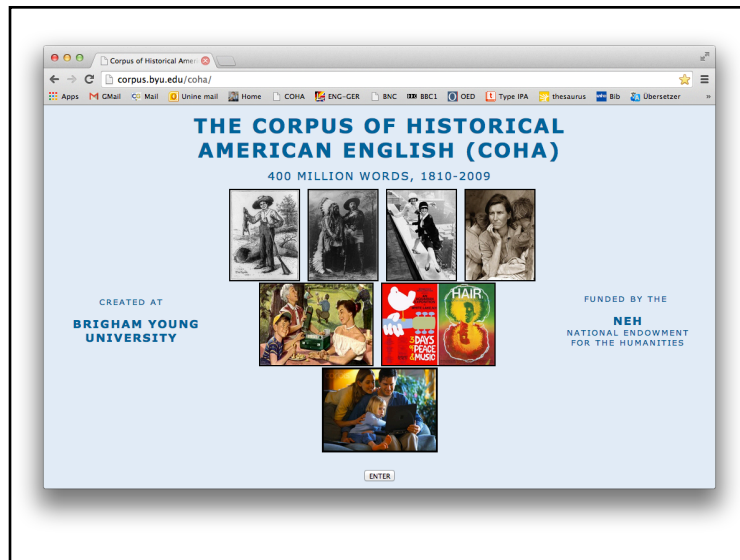


11

Questions:

How does noun-participle compounding develop? Are those developments related to changes in the English passive?

12



13

data

- examples from COHA
 - words that were tagged as adjectives
 - words that contained a hyphen
 - words that ended in a letter combination found in past participles
 - government-funded, weather-beaten, sorrow-bent, Florida-born, U.S.-built, etc.

14

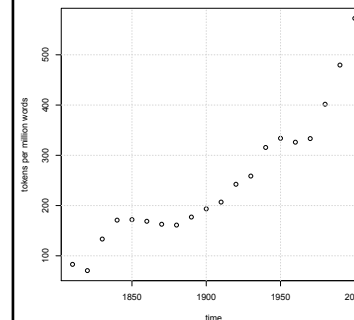
150,000 tokens

31,000 types

TYPE	1810	1820	1830	1840	1850	1860	1870	1880	1890	1900	1910	1920	1930	1940	1950	1960	1970	1980	1990	2000
GOD-ABANDONED	0	0	0	0	0	1	0	0	0	0	0	1	0	0	0	0	0	0	0	0
GOODS-ABANDONED	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
HEAVEN-ABANDONED	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
SELF-ABANDONED	0	1	0	1	0	0	1	0	0	0	0	0	0	0	1	0	0	0	0	0
SNAIL-ABANDONED	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
SOUL-ABANDONED	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
WARD-ABANDONED	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
SELF-ABASED	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
AWE-ABATED	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0
TAK-ABATED	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
WARNING-ABETTED	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
GOD-ABHORRED	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
SELF-ABOLISHED	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0
SELF-ABROGATED	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0
BOOK-ABSORBED	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0
BUSINESS-ABSORBED	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0
CAREER-ABSORBED	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0
SELF-ABSORBED	0	0	1	1	1	1	3	3	6	4	5	7	3	5	3	4	10	16	22	33
SHOCK-ABSORBED	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0
THOUGHT-ABSORBED	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0
SELF-ABSTRACTED	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0
MAN-ABUSED	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0
WAVE-ABUSED	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0
WORK-ABUSED	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0

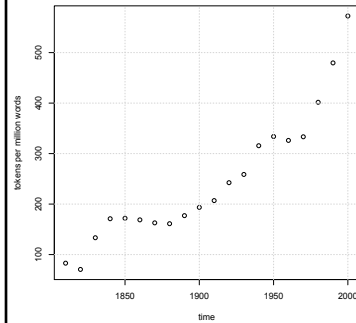
15

increase in token frequency

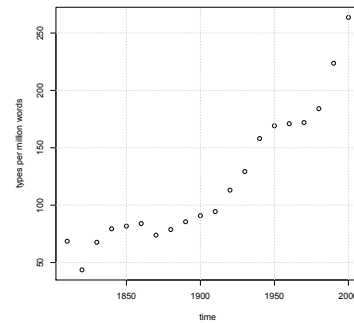


16

increase in token frequency



increase in type frequency

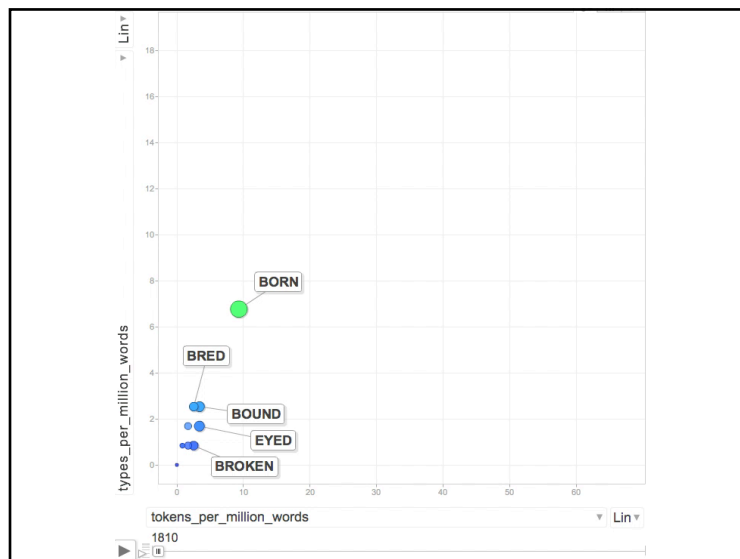


17

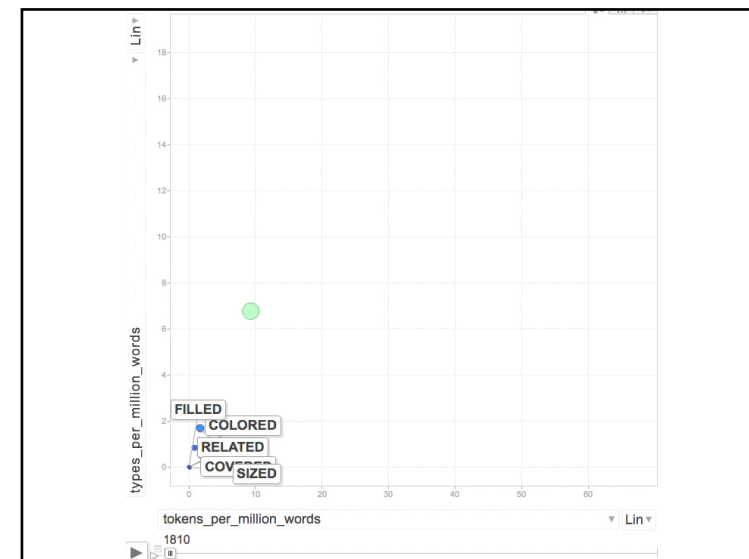
change in the participles

- How does the increase in token and type frequencies come about?
- Which are the participles that carry the increase?
- Participle families of different sizes:
 - wetted: dew, gall, snow, tear
 - yellowed: age, fear, opium, silt, smoke, sun, time, tobacco
 - coated: aluminum, bearskin, beech, blood, candy, caramel, carbon, cement, ... (129 members)

18



19



20

summarizing the developments

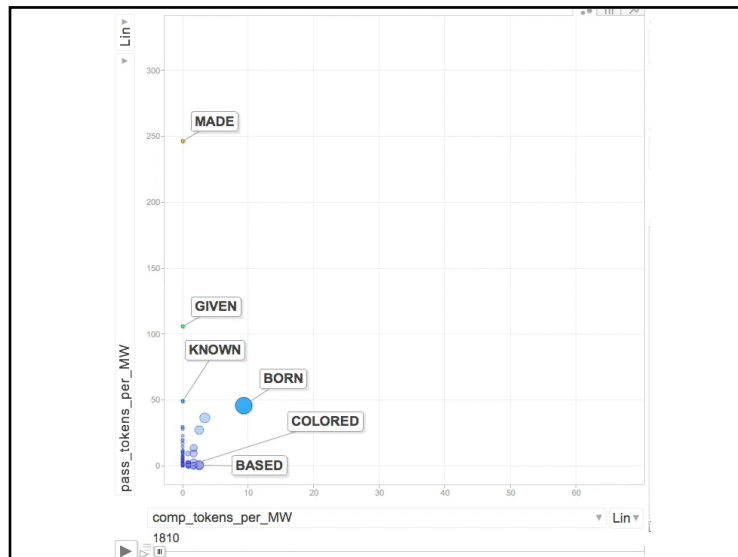
- 1810s
 - Oregon-born, Harvard-bred, context-bound, doe-eyed
- 1900s
 - honey-colored, arrow-shaped, chocolate-covered
- 2000s
 - Houston-based, work-related, toddler-sized

21

How does this compare to the passive?

- The passive with *be* in COHA:
 - [be] [v?n*]
 - 5280 types, approx. 3M tokens
- Overlapping participle types were identified:
 - government-*funded*
 - It was *funded* by the government.
- Do the overlapping types show similar frequency developments?

22



23

interpretation

- The two constructions change independently.
- The participle types that stand out most in the history of noun-participle compounding do not correspond to passive sentences.
 - *The company is based in Houston.*
 - *The problem is related to drug abuse.*
 - *The car is sized just right.*
- >> Both constructions inherit characteristics from the participle, but beyond that, speakers treat them as **two separate constructions**.

24

2. many a NOUN

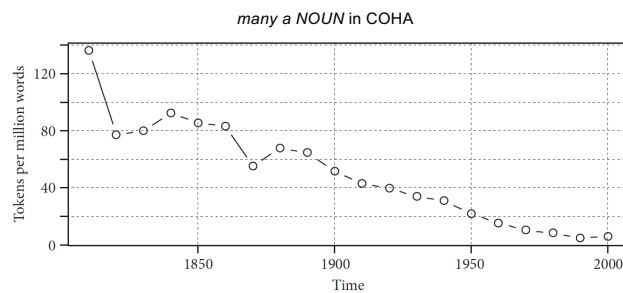
25

many a NOUN

- Time nouns
 - Many a day will pass until this construction is properly understood.
 - I've thought that many a time myself.
- Human beings
 - College education was at a premium, with many a father resisting education for a daughter.
 - Many a labour voter is not happy with the outcome.
- Totally random nouns
 - During my time in Australia I enjoyed many a sausage roll for brekkie.

26

frequencies fall, but productivity stays high



27

Question:

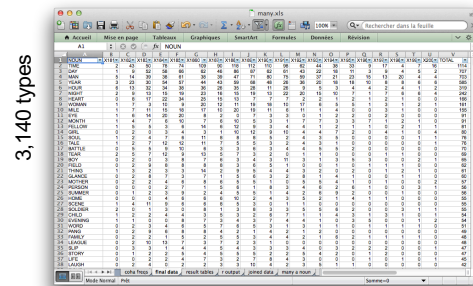
What happens to the semantic spectrum of *many a NOUN*? Why can speakers still say things like *many a sausage roll* ?

28

data

- examples from COHA
 - all sequences of *many*, *a* / *an*, and a noun
 - focus on the 230 most frequent types (63% of the data)

15,000 tokens



29

analysis

- For the 230 most frequent types, a **semantic vector space** was constructed on the basis of synchronic corpus data.
- What are the relative similarities in a group of words?

30

analysis

- For the 230 most frequent types, a **semantic vector space** was constructed on the basis of synchronic corpus data.
- What are the relative similarities in a group of words?
- Words that occur with the same collocates are judged to be similar.
- Given a word such as CHURCH, what are the lexical words that co-occur with it in a window of 4L and 4R?

31

BNC collocation frequencies (four-left, four-right window)

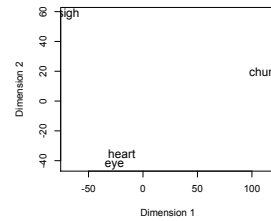
		terms			
collocates		church	heart	eye	sigh
	abbey	30	0	1	0
	above	13	8	26	0
	activities	21	2	0	0
	always	28	28	23	5
	half	8	0	10	0
	long	28	11	7	48
	family	47	10	4	1
	gave	9	13	12	109
	christ	141	8	1	0
	walk	15	2	1	0

32

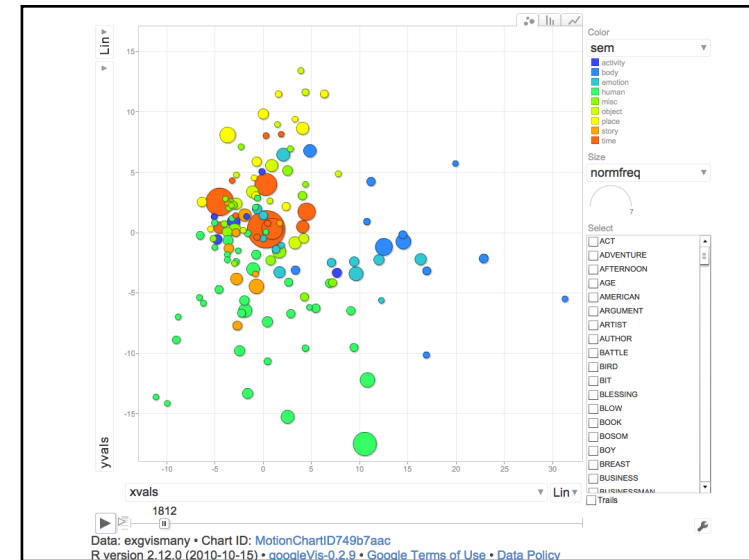
from collocate frequencies to a map

	church	heart	eye	sigh
abbey	30	0	1	0
above	13	8	26	0
activities	21	2	0	0
always	28	28	23	5
half	8	0	10	0
long	28	11	7	48
family	47	10	4	1
gave	9	13	12	109
christ	141	8	1	0
walk	15	2	1	0

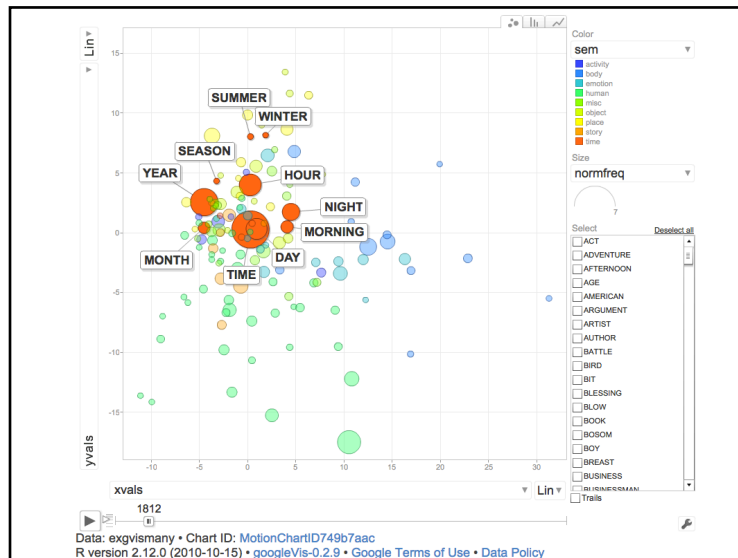
church heart eye
 heart 144.50606
 eye 153.54153 23.60085
 sigh 186.34645 106.44717 110.46266



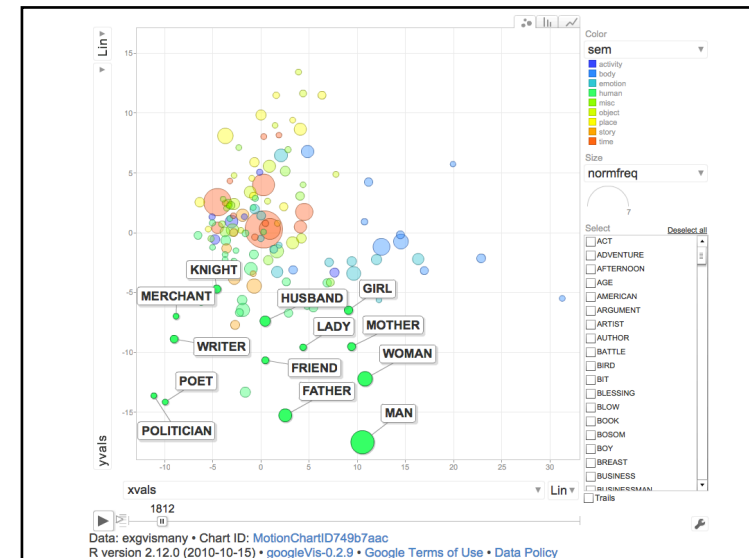
33



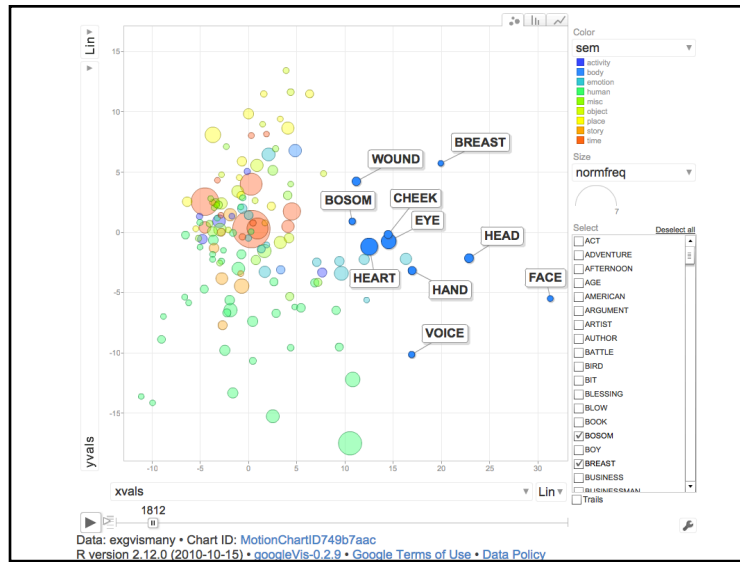
34



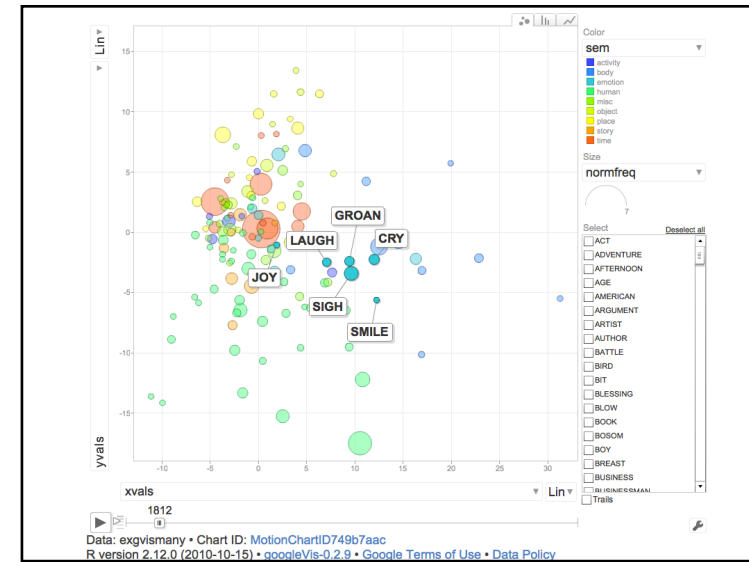
35



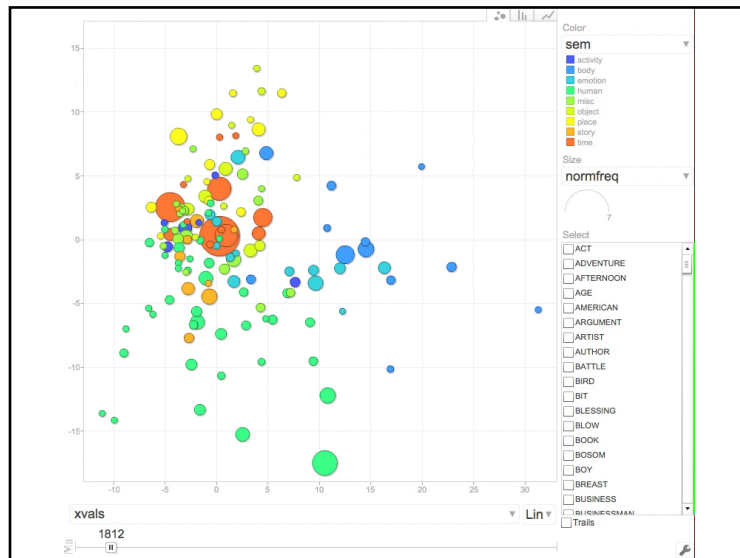
36



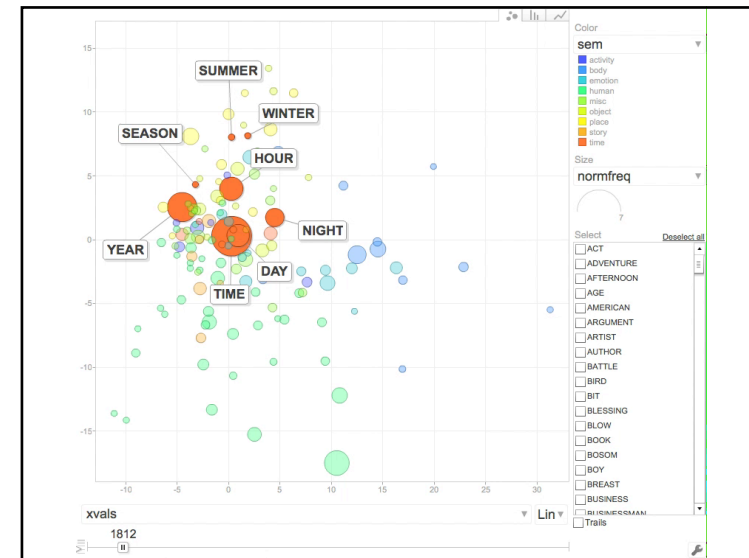
37



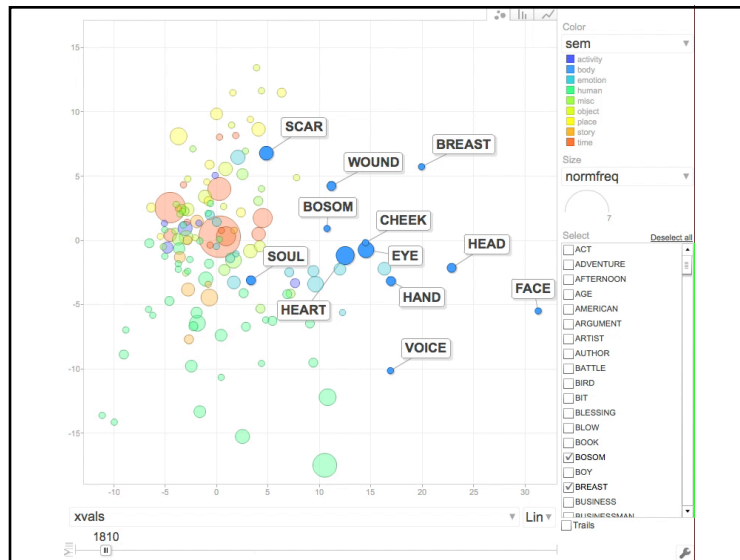
38



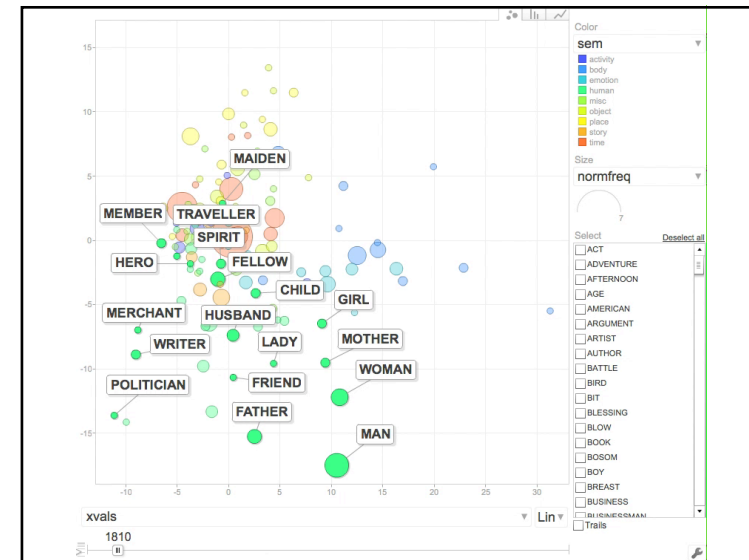
39



40



41



42

interpretation

- Whence *many a sausage roll*?
- The *many a noun* construction does not recede into a single semantic niche.
- Time nouns remain strong, human being nouns remain strong.
- Words like *time* or *man* are highly diffuse in their collocational behavior, besides them there is a sizable residue of semantically diverse types.
- Speakers thus experience the *many a NOUN* construction as **semantically unrestricted**.

43

Making your own motion charts

44

what you need

- a corpus with comparable data from different historical periods
 - Mark Davies' corpora
- the open-source software R
 - the package "googleVis"
- a spreadsheet software
- internet access



45

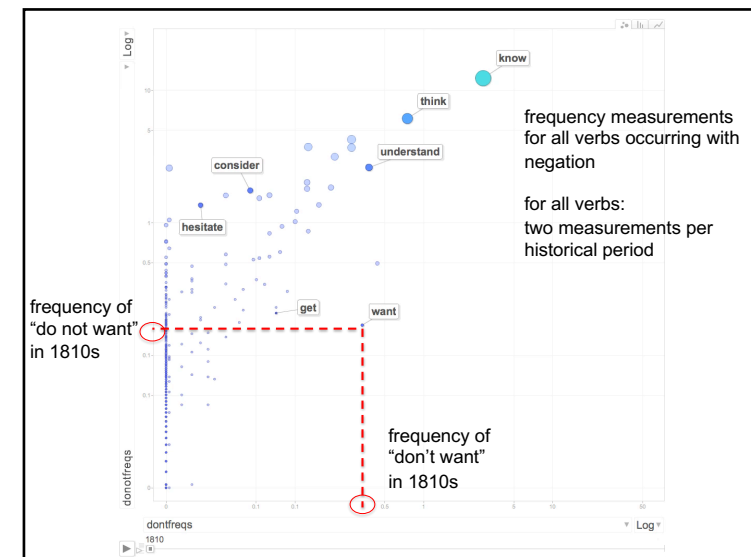
data: what kind of data, how to organize it

46

What kind of corpus data?

- easiest scenario:
 - a comparison of two alternative linguistic forms
 - negation:
 - contracted form plus verb: *don't VERB*
 - non-contracted form plus verb: *do not VERB*
 - what is compared is the frequency of elements in the VERB slot
 - Are certain verbs more likely to occur with the contracted form? How does this change over time?
 - contracted frequencies on the x-axis
 - non-contracted frequencies on the y-axis

47



48

retrieving data
do not / don't know, want, think, get, understand

CORPUS OF HISTORICAL AMERICAN ENGLISH

400 MILLION WORDS, 1810-2009

[DOWNLOAD ALL 115,000 TEXTS]

[START]

SEE CONTEXT: CLICK ON WORD (ALL SECTIONS), NUMBER (ONE SECTION), OR [CONTEXT] (SELECT)

[HELP...]

SEARCH STRING

WORD(S) [do not get/want/under]

COLLOCATES

POS LIST

[SHOW]

1 [DO] [NOT] [KNOW]

2 [DO] [NOT] [WANT]

3 [DO] [NOT] [THINK]

4 [DO] [NOT] [UNDERSTAND]

5 [DO] [NOT] [GET]

ALL

1810

1820

1830

1840

1850

1860

1870

1880

1890

1900

1910

1920

1930

1940

1950

1960

1970

1980

1990

2000

2010

2020

81259

45

139

575

734

1588

2056

2569

2840

2674

3606

4557

4957

5620

6575

7375

8125

8975

9825

10675

11525

12375

13225

14075

39005

12

14

103

190

402

622

840

1081

1451

2023

2257

2402

2402

2402

2402

2402

2402

2402

2402

2402

2402

2402

2402

25249

13

35

200

199

460

694

776

1009

1221

1296

1286

1286

1286

1286

1286

1286

1286

1286

1286

1286

1286

1286

1286

10156

6

10

65

35

100

140

136

201

194

325

560

549

735

818

901

984

1067

1150

1233

1316

1399

1482

1565

7077

4

9

34

26

85

147

217

265

311

357

426

364

395

426

457

488

519

550

581

612

643

674

705

CORPUS OF HISTORICAL AMERICAN ENGLISH

400 MILLION WORDS, 1810-2009

[DOWNLOAD ALL 115,000 TEXTS]

[START]

SEE CONTEXT: CLICK ON WORD (ALL SECTIONS), NUMBER (ONE SECTION), OR [CONTEXT] (SELECT)

[HELP...]

SEARCH STRING

WORD(S) [do not get/want/under]

COLLOCATES

POS LIST

[SHOW]

1 [DO] [NOT] [KNOW]

2 [DO] [NOT] [WANT]

3 [DO] [NOT] [THINK]

4 [DO] [NOT] [UNDERSTAND]

5 [DO] [NOT] [GET]

ALL

1810

1820

1830

1840

1850

1860

1870

1880

1890

1900

1910

1920

1930

1940

1950

1960

1970

1980

1990

2000

2010

2020

25272

32

321

633

675

1130

1296

1790

1848

1876

2215

2024

1905

1702

1584

1466

1348

1230

1112

1000

888

776

664

552

440

8367

3

36

67

101

115

154

307

347

346

479

504

609

622

635

648

661

674

687

700

713

726

739

752

7488

21

184

307

345

570

487

607

640

591

537

440

466

391

316

241

166

91

16

10

4

1

1

1

1

4324

20

107

127

136

192

224

285

283

322

363

338

267

244

221

198

175

152

129

106

84

62

40

18

1988

6

20

31

48

64

106

103

120

137

150

192

151

128

105

82

60

38

16

11

7

4

2

1

49

paste data into a spreadsheet,
bring verbs into identical sequence

negation.xlsx

Tableaux Graphiques SmartArt Formules Données Révision

Tableau 1

	1810	1820	1830	1840	1850	1860	1870	1880	1890	1900	1910	1920	1930	1940	1950	1960	1970	1980	1990	2000	2010	2020
1 [DO] [NOT] [KNOW]	45	139	575	734	1588	2056	2569	2840	2674	3606	4557	4957	5620	6575	7375	8125	8975	9825	10675	11525	12375	13225
2 [DO] [NOT] [WANT]	12	14	103	190	402	622	840	1081	1451	2023	2257	2402	2402	2402	2402	2402	2402	2402	2402	2402	2402	2402
3 [DO] [NOT] [THINK]	13	35	200	199	460	694	776	1009	1221	1296	1286	1286	1286	1286	1286	1286	1286	1286	1286	1286	1286	1286
4 [DO] [NOT] [GET]	6	10	65	35	100	140	136	201	194	325	560	549	735	818	901	984	1067	1150	1233	1316	1399	1482
5 [DO] [NOT] [UNDERSTAND]	4	9	34	26	85	147	217	265	311	357	426	464	502	540	578	616	654	692	730	768	806	844

50

getting the data into the right format
(requires some manual copy & paste)

verb	time period	contracted	non-contracted
know	1810	45	32
want	1810	12	3
think	1810	13	21
get	1810
understand	1810
know	1820	139	321
want	1820
think
...

51

problem: If the corpus periods have different sizes,
the frequencies need to be normalized!

verb	time period	contracted	non-contracted
know	1810	45	32
want	1810	12	3
think	1810	13	21
get	1810
understand	1810
know	1820	139	321
want	1820
think
...

52

period sizes of the COHA

Download all 115,000 texts, for use on your own computer.

Composition of corpus

DECADE	FICTION	POPULAR MAGAZINES	NEWSPAPERS	NON-FICTION BOOKS	TOTAL	% FICTION
1810s	641,164	88,316	0	451,542	1,181,022	0.54
1820s	1,751,204	1,714,789	0	1,461,012	6,927,005	0.54
1830s	7,590,350	3,145,575	0	3,038,062	13,773,987	0.55
1840s	8,850,886	3,554,534	0	3,641,434	16,046,854	0.55
1850s	9,094,346	4,220,558	0	3,178,922	16,493,826	0.55
1860s	9,450,562	4,437,941	262,198	2,974,401	17,125,102	0.55
1870s	10,291,968	4,452,192	1,030,560	2,835,440	18,610,160	0.55
1880s	11,215,065	4,481,568	1,355,456	3,820,766	20,872,855	0.54
1890s	11,212,219	4,679,486	1,383,948	3,907,730	21,183,383	0.53
1900s	12,029,439	5,062,650	1,433,576	4,015,567	22,541,232	0.53
1910s	11,935,701	5,694,710	1,489,942	3,534,899	22,655,252	0.53

53

45 examples in a 1.18 million word corpus:
38.1 instances per million words

verb	time period	contracted	non-contracted
know	1810	45	32
want	1810	12	3
think	1810	13	21
get	1810
understand	1810
know	1820	139	321
want	1820
think
...

139 examples in a 6.92 million word corpus:
20.6 instances per million words

54

suggested format, saved as .csv

negdata.csv

Rechercher dans la feuille

Accueil Mise en page Tableaux Graphiques SmartArt

A1 C D E F G H

verb decade full contr periodsize fullnorm contrnorm

1	know	1810	32	45	1181022	27.10	38.10
2	want	1810	3	12	1181022	2.54	10.16
3	think	1810	21	13	1181022	17.78	11.01
4	get	1810	0	6	1181022	0.00	5.08
5	understand	1810	20	4	1181022	16.93	3.39
6	know	1820	321	139	6927005	46.34	20.07
7	want	1820	36	14	6927005	5.20	2.02
8	think	1820	184	35	6927005	26.56	5.05
9	get	1820	6	10	6927005	0.87	1.44
10	understand	1820	100	9	6927005	14.44	1.30
11	know	1830	633	575	13773987	45.96	41.75
12	want	1830	67	103	13773987	4.86	7.48
13	think	1830	307	200	13773987	22.29	14.52
14	get	1830	20	65	13773987	1.45	4.72
15	understand	1830	127	34	13773987	9.22	2.47
16	know	1840	675	734	16046854	42.06	45.74
17	want	1840	81	190	16046854	5.05	11.84
18	think	1840	345	199	16046854	21.50	12.40
19	get	1840	31	35	16046854	1.93	2.18
20	understand	1840	136	26	16046854	8.48	1.62
21	know	1850	1130	1588	16493826	68.51	96.28
22	want	1850	115	402	16493826	6.97	24.37
23	think	1850	570	460	16493826	34.56	27.89
24	get	1850	48	100	16493826	2.91	6.06
25	understand	1850	192	85	16493826	11.64	5.15
26	know	1860	1296	2056	17125102	75.68	120.06
27	want	1860	154	622	17125102	8.99	36.32
28	think	1860	487	694	17125102	28.44	40.53

negdata.csv

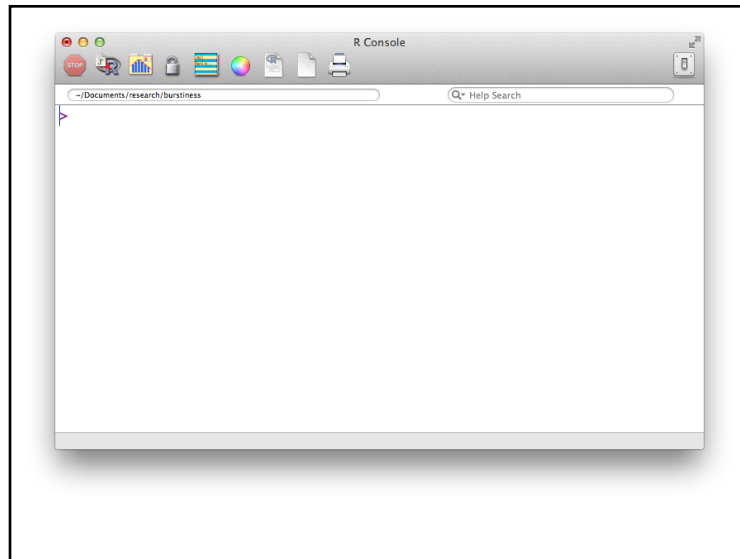
Mode Normal Print

55

creating the motion chart:
install and open R (www.r-project.org)

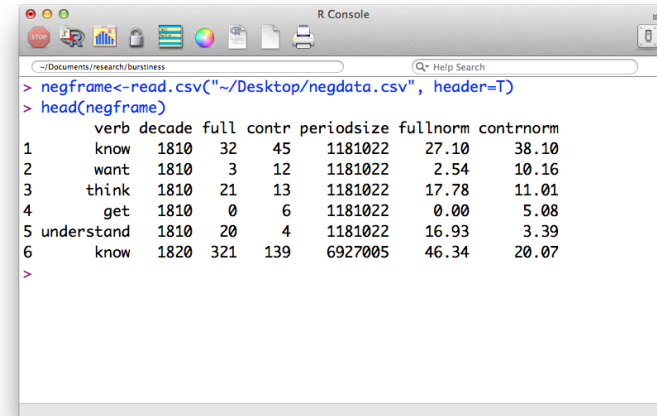


56

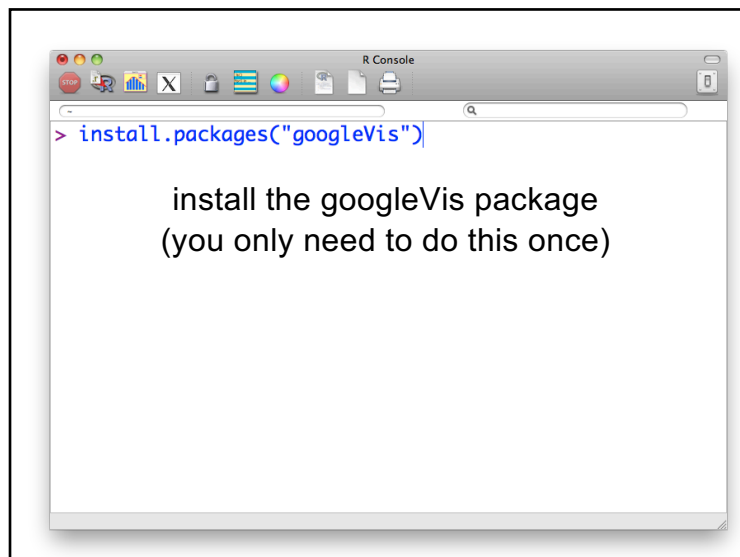


57

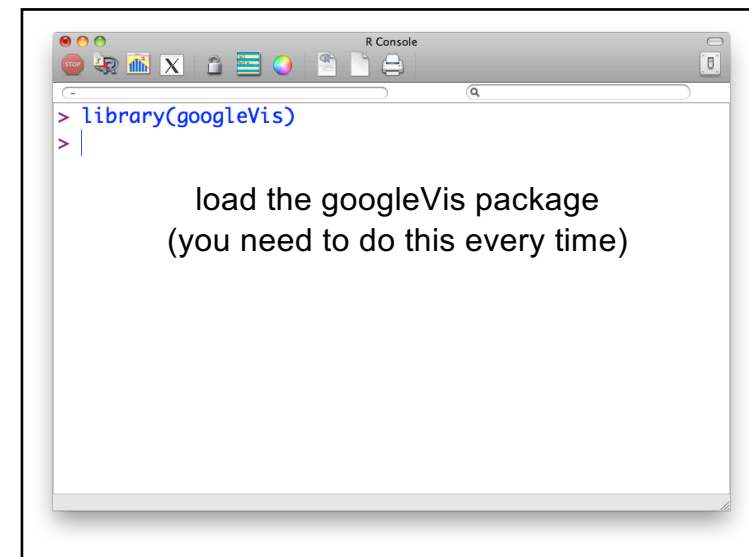
read the .csv file into R



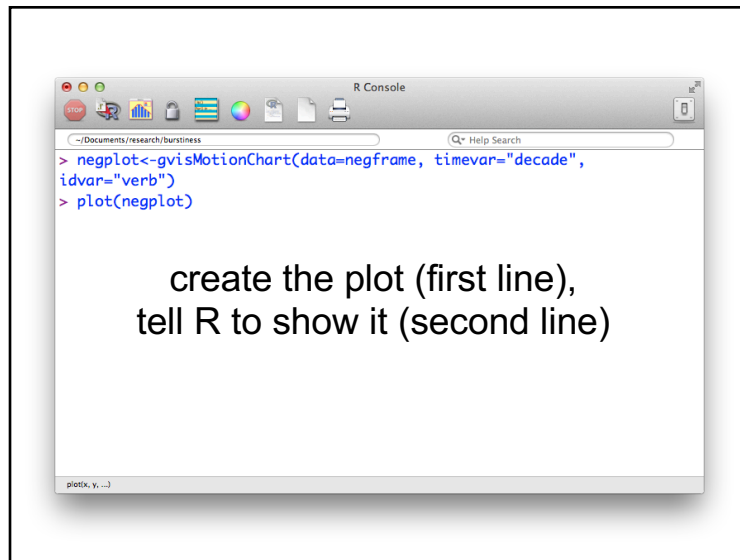
58



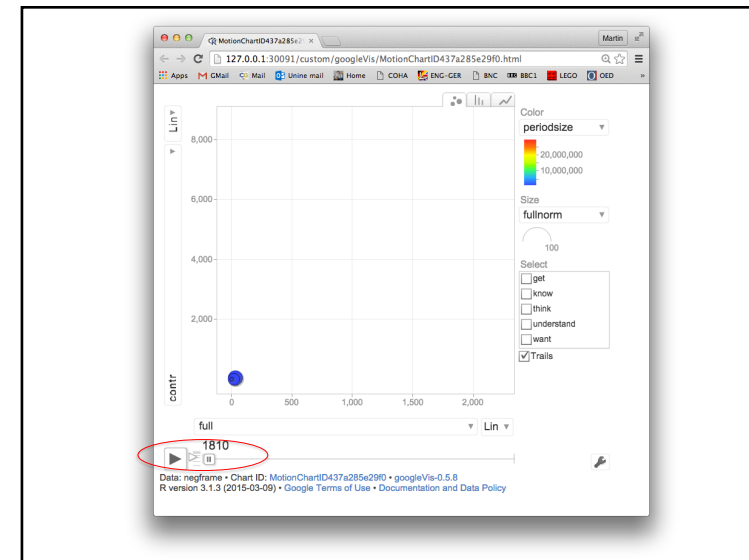
59



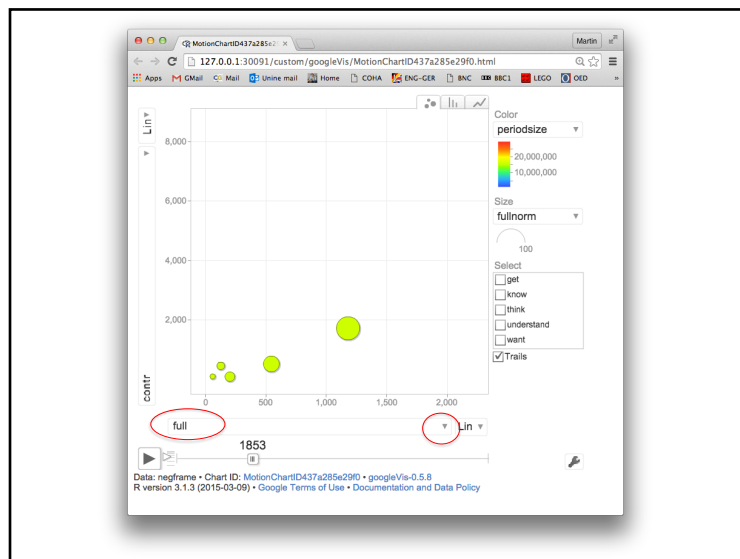
60



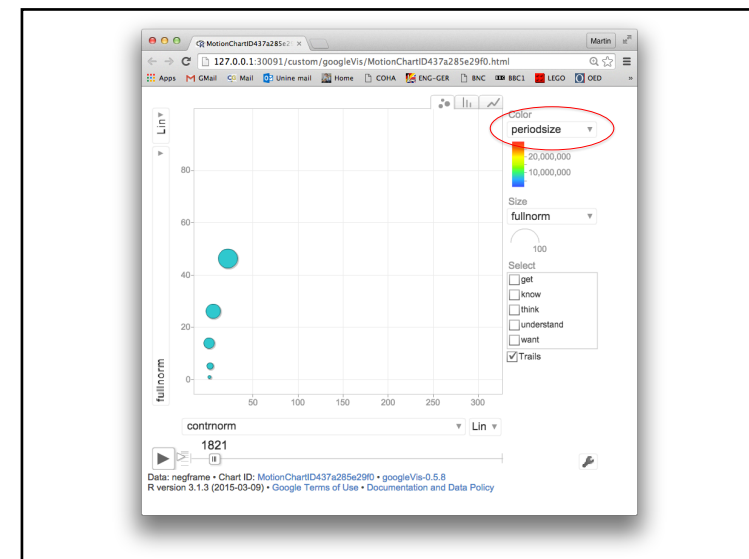
61



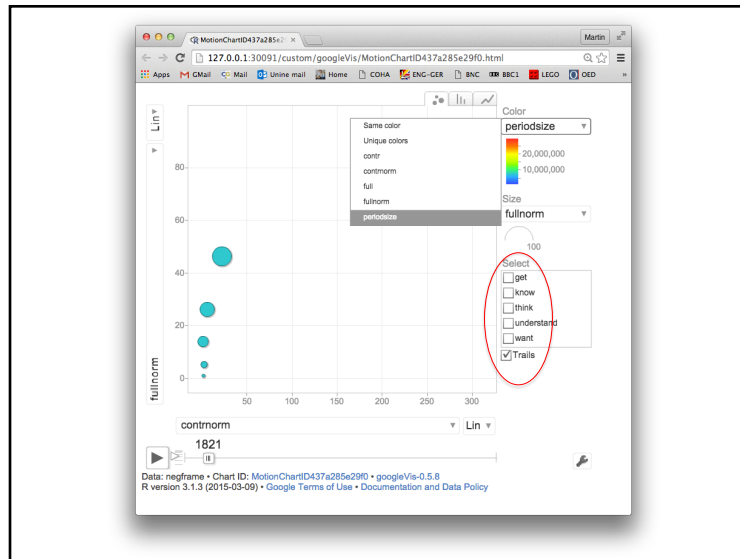
62



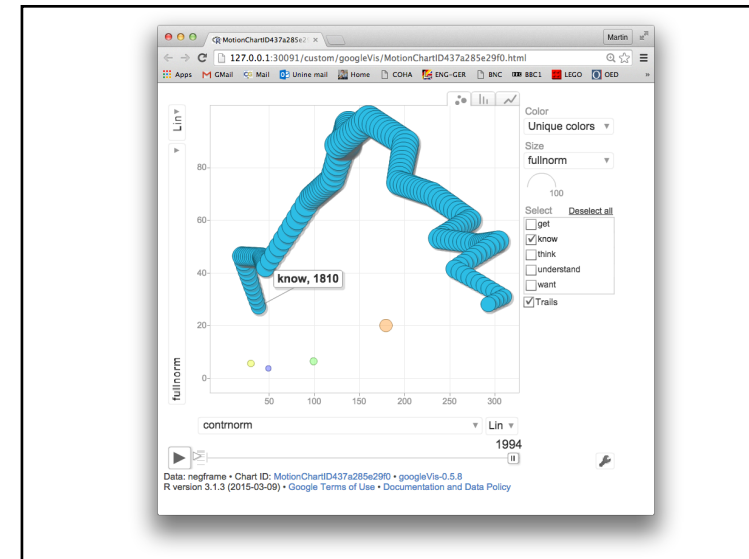
63



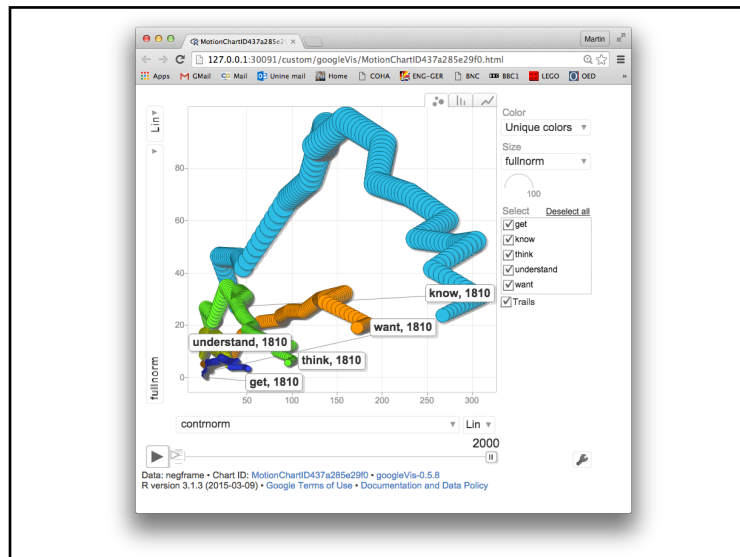
64



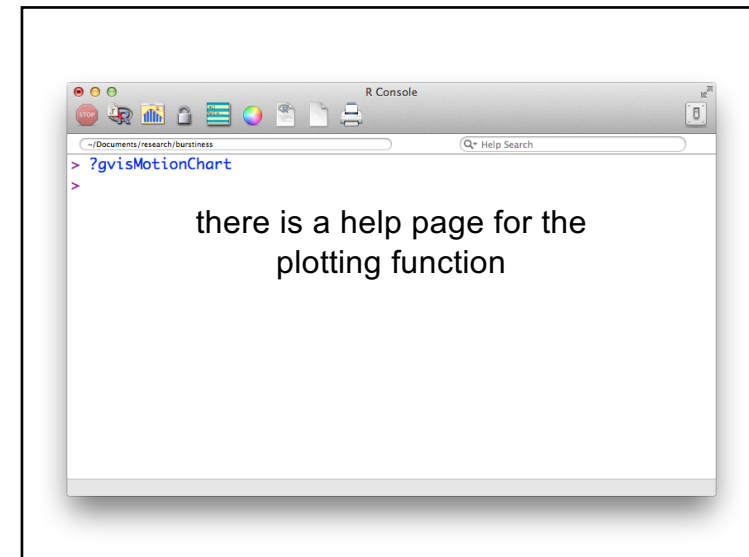
65



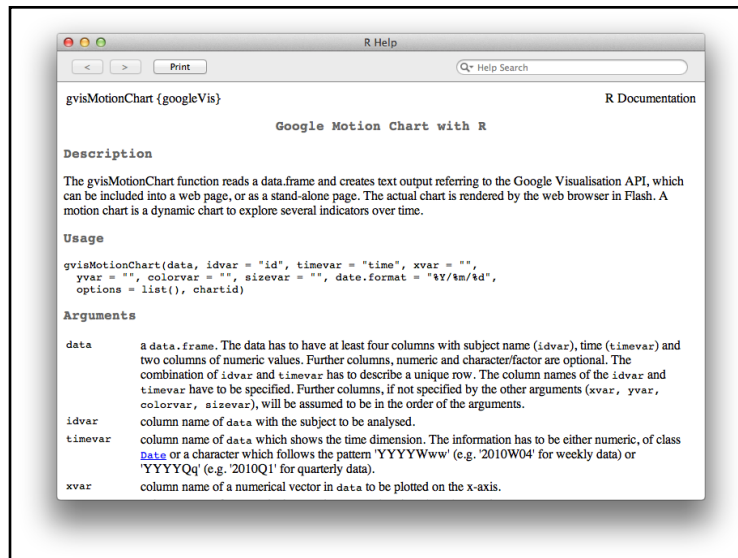
66



67



68

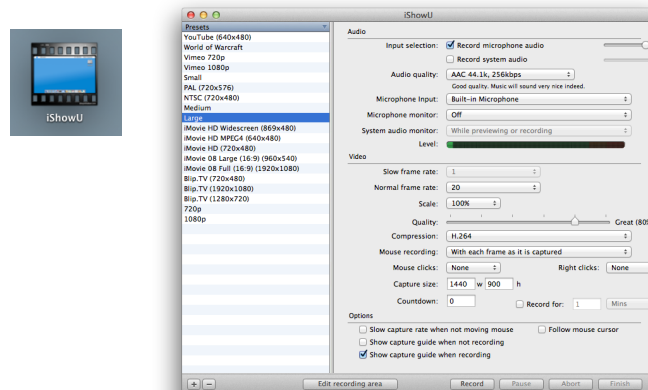


69

How can I present a motion chart when I give a talk?

70

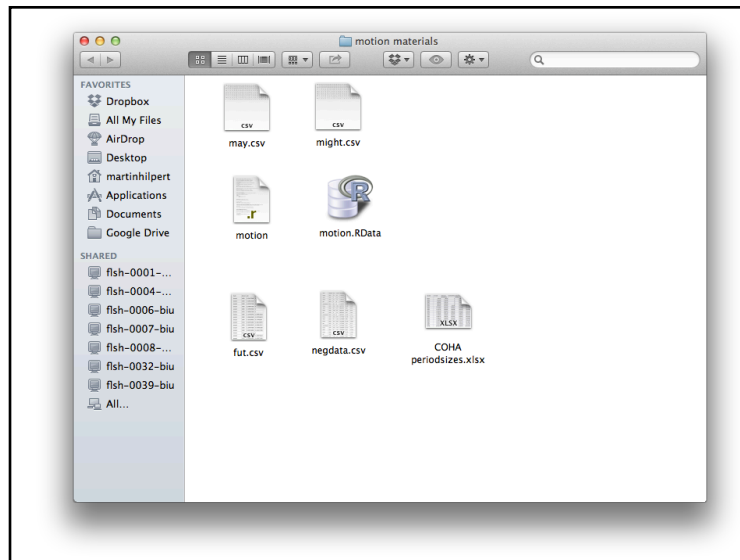
screen capture software, creates .mov files



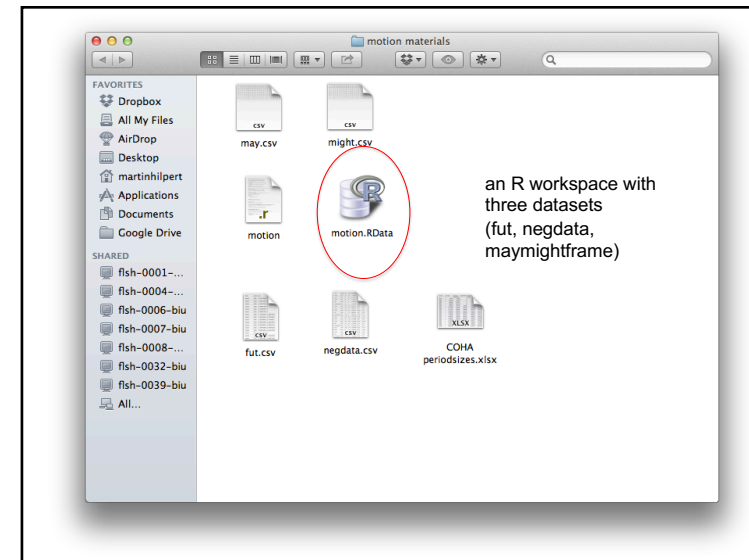
71

Some exercise materials:
<http://members.unine.ch/martin.hilpert/motion.zip>

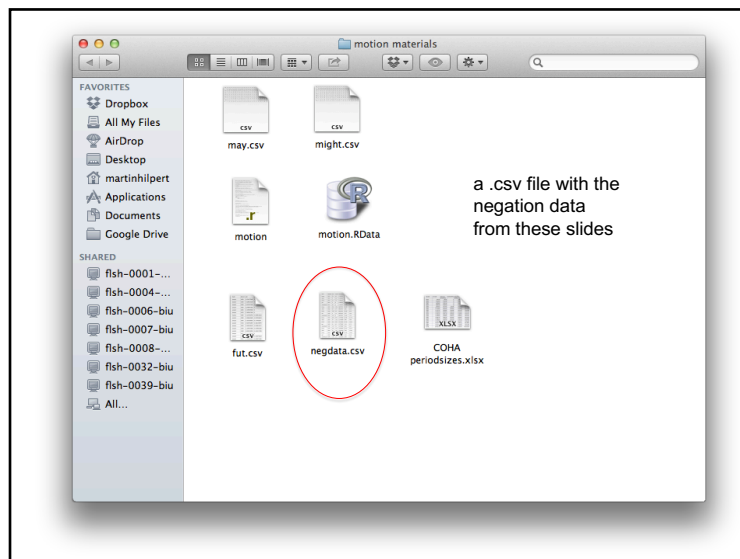
72



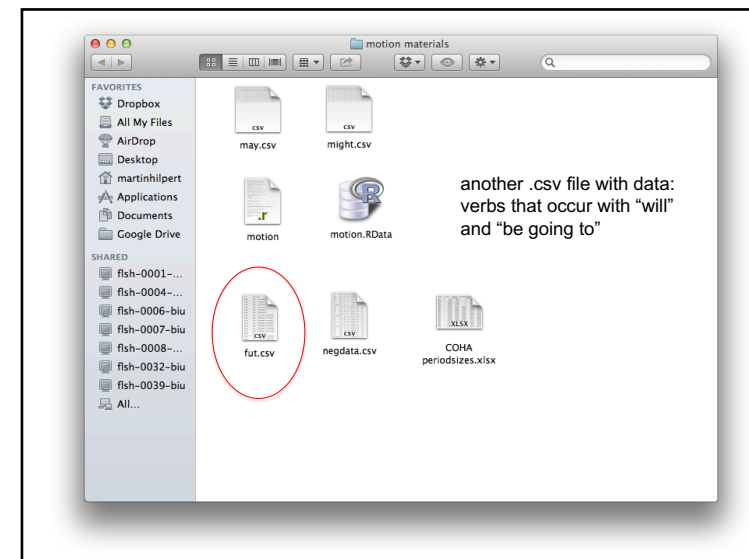
73



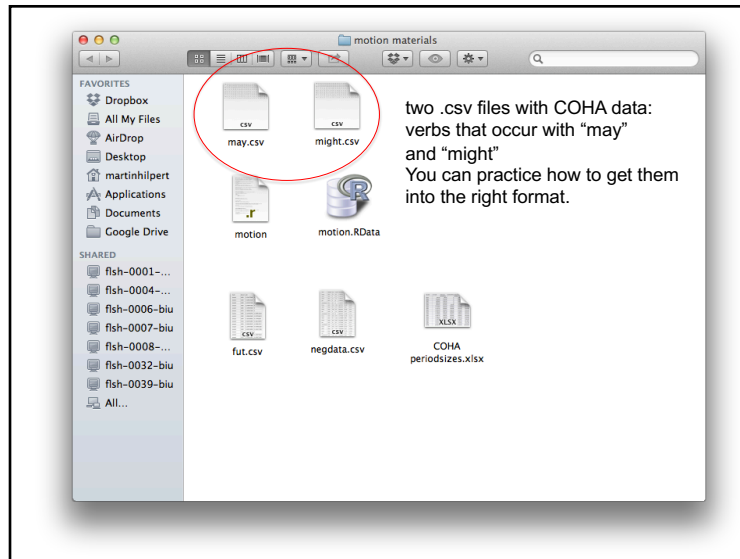
74



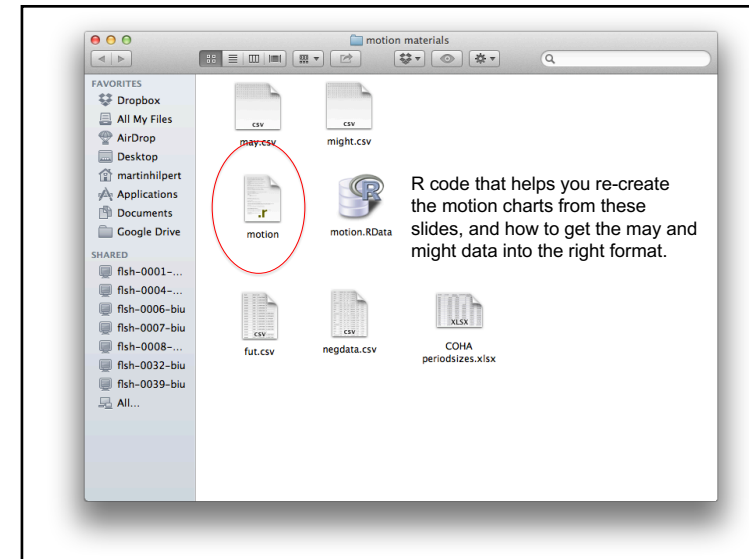
75



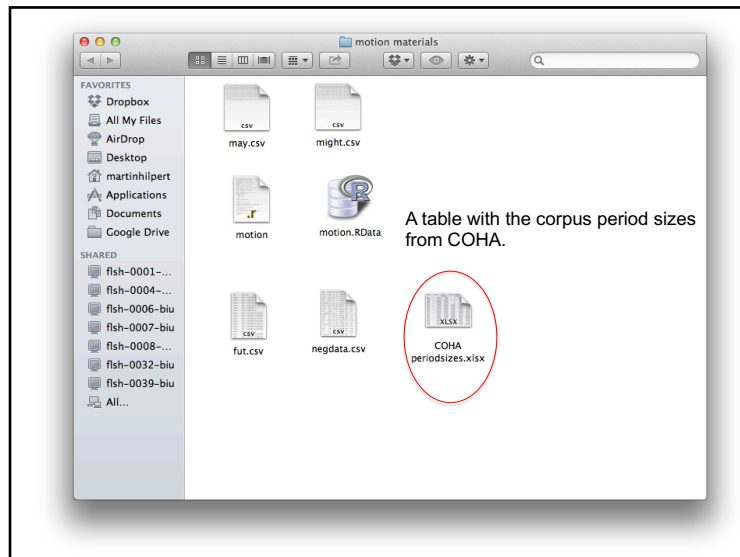
76



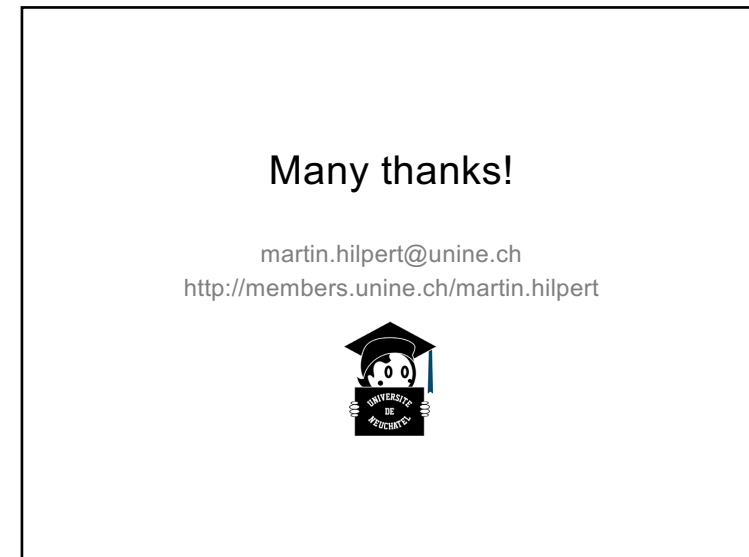
77



78



79



80