

Quantitative Discovery of Qualitative Information: A General Purpose Document Clustering Methodology

Gary King

Institute for Quantitative Social Science
Harvard University

Talk at Washington University, St. Louis, 1/21/2010

Joint work with Justin Grimmer, Harvard University

A Method for Conceptualization

- Systematic method for computer-assisted conceptualization from text

A Method for Conceptualization

- **Systematic method for computer-assisted conceptualization from text**
- Conceptualization through **Classification**: “one of the most central and generic of all our conceptual exercises. . . . the foundation not only for conceptualization, language, and speech, but also for mathematics, statistics, and data analysis. . . . Without classification, there could be no advanced conceptualization, reasoning, language, data analysis or, for that matter, social science research.” (Bailey, 1994).

A Method for Conceptualization

- **Systematic method for computer-assisted conceptualization from text**
- Conceptualization through **Classification**: “one of the most central and generic of all our conceptual exercises. . . . the foundation not only for conceptualization, language, and speech, but also for mathematics, statistics, and data analysis. . . . Without classification, there could be no advanced conceptualization, reasoning, language, data analysis or, for that matter, social science research.” (Bailey, 1994).
- We focus on **Cluster Analysis**: simultaneously 1) invent categories and 2) assign documents to categories

A Method for Conceptualization

- **Systematic method for computer-assisted conceptualization from text**
- Conceptualization through **Classification**: “one of the most central and generic of all our conceptual exercises. . . . the foundation not only for conceptualization, language, and speech, but also for mathematics, statistics, and data analysis. . . . Without classification, there could be no advanced conceptualization, reasoning, language, data analysis or, for that matter, social science research.” (Bailey, 1994).
- We focus on **Cluster Analysis**: simultaneously 1) invent categories and 2) assign documents to categories
- (We focus on texts, our methods apply more broadly)

Why Johnny Can't Classify (Optimally)

Why Johnny Can't Classify (Optimally)

- Clustering seems easy; its not!

Why Johnny Can't Classify (Optimally)

- Clustering seems easy; its not!
- $Bell(n)$ = number of ways of partitioning n objects

Why Johnny Can't Classify (Optimally)

- Clustering seems easy; its not!
- $Bell(n)$ = number of ways of partitioning n objects
- $Bell(2) = 2$ (AB, A B)

Why Johnny Can't Classify (Optimally)

- Clustering seems easy; its not!
- $Bell(n)$ = number of ways of partitioning n objects
- $Bell(2) = 2$ (AB, A B)
- $Bell(3) = 5$ (ABC, AB C, A BC, AC B, A B C)

Why Johnny Can't Classify (Optimally)

- Clustering seems easy; its not!
- $Bell(n)$ = number of ways of partitioning n objects
- $Bell(2) = 2$ (AB, A B)
- $Bell(3) = 5$ (ABC, AB C, A BC, AC B, A B C)
- $Bell(5) = 52$

Why Johnny Can't Classify (Optimally)

- Clustering seems easy; its not!
- $Bell(n)$ = number of ways of partitioning n objects
- $Bell(2) = 2$ (AB, A B)
- $Bell(3) = 5$ (ABC, AB C, A BC, AC B, A B C)
- $Bell(5) = 52$
- $Bell(100) \approx$

Why Johnny Can't Classify (Optimally)

- Clustering seems easy; its not!
- $Bell(n)$ = number of ways of partitioning n objects
- $Bell(2) = 2$ (AB, A B)
- $Bell(3) = 5$ (ABC, AB C, A BC, AC B, A B C)
- $Bell(5) = 52$
- $Bell(100) \approx 10^{28} \times$ Number of elementary particles in the universe

Why Johnny Can't Classify (Optimally)

- Clustering seems easy; its not!
- $Bell(n)$ = number of ways of partitioning n objects
- $Bell(2) = 2$ (AB, A B)
- $Bell(3) = 5$ (ABC, AB C, A BC, AC B, A B C)
- $Bell(5) = 52$
- $Bell(100) \approx 10^{28} \times$ Number of elementary particles in the universe
- Now imagine choosing the *optimal* classification scheme by hand!

Why Johnny Can't Classify (Optimally)

- Clustering seems easy; its not!
- $Bell(n)$ = number of ways of partitioning n objects
- $Bell(2) = 2$ (AB, A B)
- $Bell(3) = 5$ (ABC, AB C, A BC, AC B, A B C)
- $Bell(5) = 52$
- $Bell(100) \approx 10^{28} \times$ Number of elementary particles in the universe
- Now imagine choosing the *optimal* classification scheme by hand!
- Its no surprise that automated algorithms can help, but which algorithms?

Why HAL Can't Classify Either

Why HAL Can't Classify Either

- Large quantitative literature on **cluster analysis**

Why HAL Can't Classify Either

- Large quantitative literature on **cluster analysis**
- The Goal — an optimal application-independent cluster analysis method — is mathematically impossible:

Why HAL Can't Classify Either

- Large quantitative literature on **cluster analysis**
- The Goal — an optimal application-independent cluster analysis method — is mathematically impossible:
 - **No free lunch theorem**: every possible clustering method performs equally well on average over all possible substantive applications

Why HAL Can't Classify Either

- Large quantitative literature on **cluster analysis**
- The Goal — an optimal application-independent cluster analysis method — is mathematically impossible:
 - **No free lunch theorem**: every possible clustering method performs equally well on average over all possible substantive applications
- Existing methods:

Why HAL Can't Classify Either

- Large quantitative literature on **cluster analysis**
- The Goal — an optimal application-independent cluster analysis method — is mathematically impossible:
 - **No free lunch theorem**: every possible clustering method performs equally well on average over all possible substantive applications
- Existing methods:
 - **Many choices**: model-based, subspace, spectral, grid-based, graph-based, fuzzy k -modes, affinity propagation, self-organizing maps,...

Why HAL Can't Classify Either

- Large quantitative literature on **cluster analysis**
- The Goal — an optimal application-independent cluster analysis method — is mathematically impossible:
 - **No free lunch theorem**: every possible clustering method performs equally well on average over all possible substantive applications
- Existing methods:
 - **Many choices**: model-based, subspace, spectral, grid-based, graph-based, fuzzy k -modes, affinity propagation, self-organizing maps, . . .
 - **Well-defined** statistical, data analytic, or machine learning foundations

Why HAL Can't Classify Either

- Large quantitative literature on **cluster analysis**
- The Goal — an optimal application-independent cluster analysis method — is mathematically impossible:
 - **No free lunch theorem**: every possible clustering method performs equally well on average over all possible substantive applications
- Existing methods:
 - **Many choices**: model-based, subspace, spectral, grid-based, graph-based, fuzzy k -modes, affinity propagation, self-organizing maps, . . .
 - **Well-defined** statistical, data analytic, or machine learning foundations
 - How to add substantive knowledge: With few exceptions, **unclear**

Why HAL Can't Classify Either

- Large quantitative literature on **cluster analysis**
- The Goal — an optimal application-independent cluster analysis method — is mathematically impossible:
 - **No free lunch theorem**: every possible clustering method performs equally well on average over all possible substantive applications
- Existing methods:
 - **Many choices**: model-based, subspace, spectral, grid-based, graph-based, fuzzy k -modes, affinity propagation, self-organizing maps, . . .
 - **Well-defined** statistical, data analytic, or machine learning foundations
 - How to add substantive knowledge: With few exceptions, **unclear**
 - The literature: **little guidance on when methods apply**

Why HAL Can't Classify Either

- Large quantitative literature on **cluster analysis**
- The Goal — an optimal application-independent cluster analysis method — is mathematically impossible:
 - **No free lunch theorem**: every possible clustering method performs equally well on average over all possible substantive applications
- Existing methods:
 - **Many choices**: model-based, subspace, spectral, grid-based, graph-based, fuzzy k -modes, affinity propagation, self-organizing maps, . . .
 - **Well-defined** statistical, data analytic, or machine learning foundations
 - How to add substantive knowledge: With few exceptions, **unclear**
 - The literature: **little guidance on when methods apply**
 - **Deriving such guidance**: difficult or impossible

Why HAL Can't Classify Either

- Large quantitative literature on **cluster analysis**
- The Goal — an optimal application-independent cluster analysis method — is mathematically impossible:
 - **No free lunch theorem**: every possible clustering method performs equally well on average over all possible substantive applications
- Existing methods:
 - **Many choices**: model-based, subspace, spectral, grid-based, graph-based, fuzzy k -modes, affinity propagation, self-organizing maps, . . .
 - **Well-defined** statistical, data analytic, or machine learning foundations
 - How to add substantive knowledge: With few exceptions, **unclear**
 - The literature: **little guidance on when methods apply**
 - **Deriving such guidance**: difficult or impossible
- **Deep problem in cluster analysis literature**: no way to know which method will work *ex ante*

If Ex Ante doesn't work, try Ex Post

If Ex Ante doesn't work, try Ex Post

- **Methods and substance must be connected** (no free lunch theorem)

If Ex Ante doesn't work, try Ex Post

- **Methods and substance must be connected** (no free lunch theorem)
- The usual approach fails: hard to do it by understanding the model

If Ex Ante doesn't work, try Ex Post

- **Methods and substance must be connected** (no free lunch theorem)
- The usual approach fails: hard to do it by understanding the model
- We do it **ex post** (by qualitative choice). For example:

If Ex Ante doesn't work, try Ex Post

- **Methods and substance must be connected** (no free lunch theorem)
- The usual approach fails: hard to do it by understanding the model
- We do it **ex post** (by qualitative choice). For example:
 - Create long list of clusterings; choose the best

If Ex Ante doesn't work, try Ex Post

- **Methods and substance must be connected** (no free lunch theorem)
- The usual approach fails: hard to do it by understanding the model
- We do it **ex post** (by qualitative choice). For example:
 - Create long list of clusterings; choose the best
 - Too hard for mere humans!

If Ex Ante doesn't work, try Ex Post

- **Methods and substance must be connected** (no free lunch theorem)
- The usual approach fails: hard to do it by understanding the model
- We do it **ex post** (by qualitative choice). For example:
 - Create long list of clusterings; choose the best
 - Too hard for mere humans!
 - An **organized** list will make the search possible

If Ex Ante doesn't work, try Ex Post

- **Methods and substance must be connected** (no free lunch theorem)
- The usual approach fails: hard to do it by understanding the model
- We do it **ex post** (by qualitative choice). For example:
 - Create long list of clusterings; choose the best
 - Too hard for mere humans!
 - An **organized** list will make the search possible
 - E.g.,: consider two clusterings that differ only because one document (of many) moves from category 5 to 6

Our Idea: Meaning Through Geography

Set of clusterings

Our Idea: Meaning Through Geography

Set of clusterings \approx

A list of unconnected addresses

wide at SuperPages.com

	195	Car	C
Cartage New England Inc 28 Allen Ln Ipswich 01938.....	978 356-9960		
Cartagena Lydia 28 Sweet Box 02331.....	617 323-7639		
Cartagena Avish F Pleasant Rd 02139.....	617 442-9780		
B Hrd 02134.....	617 361-5253		
Justica 50 Decatur Cha 02129.....	617 241-0152		
Luzilla 124 Harvard Cam 02138.....	617 491-5621		
M 95 Howe Box 02132.....	617 323-9713		
Melvin 503 Green Cam 02139.....	617 576-1061		
361-0380 Carte Nicholas 18 Appleton Boston 02114.....	617 695-6996		
Cartagena D 4 Bradford Box 02138.....	617 338-9219		
628-8248 Carten Thos Jr Sr & Claire 1 Franklin St Mt 02136.....	617 698-6163		
445-5116 Thomas & Kathleen 50 Thompson Ln Mt 02136.....	617 696-6919		
822-2962 Carte A Box 02133.....	617 229-2257		
427-5712 A Nelson A 21 Bethune Wy Roxbury 02119.....	617 442-1219		
569-2698 A 200 Pines Av Cambridge 02142.....	617 492-4174		
667-5190 A M 255 Massachusetts Av Box 02115.....	617 266-7153		
569-1412 Adams 361 Centre St Mt 02136.....	617 698-7074		
338-9110 Alice 108 Elmwood Box 02114.....	617 423-0193		
825-9119 Andrew F 42 West St Box 02143.....	617 945-2711		
825-9119 Carte Anne MD 1161 Beacon St 02444.....	617 739-1022		
296-1593 Carte J M 371 Newbury Boston 02116.....	617 536-6329		
670-2078 B E 10 Graduate Av Mt 02136.....	617 296-6911		
621-9001 Carte Barbara L MD Tufts New England Medical Center Box 02111 Cam.....	617 636-0051		
296-4725 Carte Becky Box 02114.....	617 523-4368		
542-1521 Bernard J 122 Goodhue F Rd 02136.....	617 567-9430		
364-5232 Bibbiah 25 Midway Dr 02134.....	617 298-8713		
541-5429 Bilal 26 Elmwood St 02138.....	617 367-9931		
739-2662 Carte Broadcasting Co 58 Park Pl Box 02114.....	617 423-0210		
879-0030 Carte C 2000 Commonwealth Av 02135.....	617 225-0200		
436-1511 C 218 Harvard Av East Boston 02128.....	617 569-1545		
569-4119 C 109 Harvard Cam 02138.....	617 491-4822		
809-6212 C 28 Irving St Cambridge 02142.....	617 526-4392		
869-8782 C & M 43 Bernham Jan 02136.....	617 524-9558		
317 327-1105 Carter F 24 Hibisc Box 02131.....	617 327-1105		
Faye & Ricky 27 Columbia Av Box 02136.....	617 437-7331		
Francis S 134 Temple W Av 02132.....	617 323-6781		
Franklin & Anne 75 Mt Auburn Cam 02138.....	617 354-0798		
Fred 42 Harvard Jan 02136.....	617 524-3078		
Fred 16 Newbury Av Mt 02136.....	617 698-1343		
G & B 8 Verden Dr 02134.....	617 434-8966		
G T 27 Franklin Av Sun 02145.....	617 623-7121		
Gayle 25 Franklin Dr 02134.....	617 823-0322		
Geo S 115 Mass Mt Jan 02138.....	617 522-3215		
George 125 Madison Box 02114.....	617 367-9548		
Carter Hillside Assoc 107 S Street Box 02111.....	617 456-1689		
Carter Harry F 36 Bayview Rd W Av 02132.....	617 325-5465		
Carter Hide Co Inc 100 Riverside Dr 02148.....	617 542-7987		
Carter Hilary 41 Harvey Cam 02148.....	617 876-2750		
Horace 361 Walnut Av Roxbury 02119.....	617 442-5307		
Howard Jr 28 Nona Dr Box 02118.....	617 445-5532		
J Dan 41 Chatham Box 02444.....	617 232-7990		
J S 538 Harvard Box 02444.....	617 730-9483		
J 775 The Pines West Roxbury 02132.....	617 323-5374		
Carter J Jacques MD 1 Ipswich Pl Box 02444.....	617 735-8787		
Carter J M 3410 Columbia Rd S Box 02137.....	617 464-1040		
Carter J M Ornamental Ironworks 200 Franklin Falls 02146.....	617 434-5353		
Carter J Veal Co 40 Newbury Rd 02138.....	617 442-1775		
Carter James 157 Cambridge St Cam 02136.....	617 492-1214		
James 412 Foster Av Roxbury 02132.....	617 739-2193		
Janet 41 Good Star Rd Cambridge 02141.....	617 876-8841		
Jane L 34 Rosbury Rd Mt 02134.....	617 361-0773		
Janice 134 Adams Rd Newton 02465.....	617 564-0435		
Jeffrey 41 Warren Av Box 02114.....	617 424-5994		
John 111 Mansfield Rd 02134.....	617 987-2163		
John 107 Summer Box 02135.....	617 423-4334		
John 40 Harvard Box 02132.....	617 282-1235		
Jane O 129 A Summit Av Box 02133.....	617 734-6109		
J 28 Irving St Cambridge 02142.....	617 265-8456		
J P 129 Harvard Cambridge 02132.....	617 282-1593		
317 267-6483 Carter Nella E 323 Main St Box 02115.....	617 267-6483		
Nicholas S F 115 Randolph Av Mt 02136.....	617 698-6307		
Nick 21 Fyfe Hill Box 02114.....	617 267-5222		
Nick & Debbi 196 Vermont Rd Newton 02459.....	617 527-0480		
Nicole 38 Chickadee Dr 02125.....	617 822-1201		
Norman G P 40 Cranston Pl Box 02135.....	617 437-4754		
P E 501 E South S Box 02137.....	617 268-8213		
P L 44 Hutchings Box 02131.....	617 427-9170		
P R 91 Boyer Jan 02134.....	617 968-8692		
Paul & Constance 114 Franklin Av W Box 02132.....	617 325-2036		
Paul E 501 E South S Box 02137.....	617 268-4546		
Paul M 27 Union Rd 02135.....	617 787-2115		
Carter Pike Driving Inc 17 Avenue Ct Framingham 02702.....	Wellesley Falls 781.235-0488		
Carter Prudence 40 Franklin Waterbury 02172.....	617 393-3782		
Prudence 40 Franklin Waterbury 02172.....	617 926-7063		
Roginald 106 Brookside Dorchester 02122.....	617 541-2843		
Renee & Andrew 10 Walnut Box 02138.....	617 720-3765		
Carter Rice David Bulfinch Boston Publishing 163 Main Wilmington 01887 Toll Free-Dial '9 & Then.....	800 638-1671		
Carl Eric Industrial Prod 113 Main Wilmington Toll Free-Dial '9 & Then.....	800 616-7447		
Carl Toll Free-Dial '9 & Then.....	800 648-7447		
Headquarters 113 Main Wilmington 01887 Cam.....	978 988-7447		
Ingala One 163 Main Wilmington 01887 Cam.....	800 638-1673		
Carter Richard 2079 Commonwealth Av Brighton 02111.....	617 987-0836		
Richard A 97 Mt Vernon Box 02106.....	617 566-7293		
Carter Richard A MD 120 Commonwealth Pl Mt 02134.....	617 267-0710		
Carter Richard K 23 Mather S Box 02137.....	617 268-0448		
Richard L 175 Rockdale Av Cam 02141.....	617 864-1535		
Roger 130 St Braughn Box 02111.....	617 424-6148		
Roy 41 Concord Cam 02138.....	617 491-6115		
Royce 18 Saffery Cha 02129.....	617 241-0418		

Our Idea: Meaning Through Geography

Set of clusterings \approx

A list of unconnected addresses

wide at [SuperPages.com](#)

195

Car

C

17 566-1282	Cartage New England Inc 28 Allen St Ipswich 01938	978 356-9960
81 447-4101	Cartagena Lydia 28 Sweet Briar 02331	617 323-7639
90 257-9961	Cartagena Avish F Harvard St 02119	617 442-9780
	B Had 02134	617 361-5253
17 566-1282	Jesticca 50 Decatur Cha 02129	617 241-0152
17 364-5188	Luzmila 124 Harvard Cam 02136	617 491-5621
	M 95 Howe St 02136	617 323-9713
361-0380	Melvin 503 Green Cam 02139	617 576-1061
17 566-4548	Carte Nicholas 18 Appleton Boston 02114	617 695-6996
	Cartagena O 4 Bradford Bay 02118	617 338-9219
17 628-8248	Carten Thos J Sr & Claire 1 Furlow St Mt 02136	617 698-6163
17 445-5116	Thomas & Kathleen 50 Thompson Ln Mt 02136	617 696-6919
17 822-2962	Carter A An 02113	617 229-2257
17 427-5712	A Helen 116 1/2 Bedford Wy Roxbury 02119	617 442-1219
17 569-2698	A 200 Riverside Av Cambridge 02142	617 492-4174
17 667-5190	A M 255 Massachusetts Av 02115	617 266-7153
	Adams 301 Carter St Mt 02136	617 698-9074
17 569-1412	Alice 138 Elmwood Av 02118	617 453-0193
17 338-9110	Allice 40 Market Cambridge 02139	617 945-2711
	Andrew F 42 West St 02135	617 625-7623
17 825-9195	Carter Anne MD 1101 Beacon Bk 02444	617 739-1022
17 296-1593	Carter Athene 771 Newbury Boston 02116	617 536-6229
17 670-2078	B E 18 Gladstone Av Mt 02136	617 296-6911
17 621-9001	Carter Barbara L MD Tufts-New England Medical Center 02111	
17 296-4725	Call Carter Becky 02114	617 636-0951
	Carter Adhena 175 Cambridge St Cam 02138	617 523-4368
17 542-1521	Bernard J James Mathews R 02136	617 567-9430
17 364-5232	Bibbath 25 Midway Der 02124	617 298-8713
17 541-5649	Bibbath 30 New Vernon 02136	617 367-9931
17 739-2662	Carter Broadcasting Co 58 Park Pl 02116	617 423-0210
	Carter Brooks Consultants Inc 73 East C St 02471	617 225-2020
17 879-0030	Carter C 2000 Massachusetts Av 02135	617 782-2118
17 541-3948	C 210 Fenwick Av East Boston 02128	617 569-1545
17 436-1511	C 109 Harvard Cam 02138	617 491-4822
17 569-6119	C 281 Fenwick Av East Boston 02128	617 569-4392
809 669-8782	C & M 41 Northgate Jan 02124	617 524-9558
	Carter F 54 Hibiscus Bay 02131	617 327-1105
	Faye & Ricky 20 Columbia Av 02136	617 437-7331
	Francis S 134 Temple W An 02132	617 323-6781
	Franklin & Anne 751 Mt Auburn Cam 02138	617 354-0798
	Fred 41 Hawthorn Jan 02136	617 524-3078
	Fred 16 Howland Av Mt 02136	617 698-1343
	G & B 8 Verdun Der 02124	617 436-8906
	G T 27 Franklin Av Sun 02145	617 623-7121
	Gayle 25 Franklin Der 02124	617 825-8322
	Geo S 115 Mass Hill Jan 02138	617 522-3215
	George 125 Boston Bay 02134	617 367-9548
	Carter Hillside Assoc 107 S Street Bay 02111	617 456-1689
	Carter Harry F 168 Bayview Rd W An 02132	617 325-5465
	Carter Hide Co Inc 167 Essex St 02148	617 542-7987
	Carter Hilary 41 Harvey Cam 02148	617 876-2750
	Horace 301 Walnut Av Roxbury 02119	617 442-5307
	Howard Jr 28 New One Box 02118	617 445-5552
	J Can 15 Chatham Bk 02444	617 232-7990
	J 538 Harvard St 02138	617 730-9483
	J 775 The Pine Wy Roxbury 02132	617 323-5574
	Carter J Jacques MD 1 Crockett Pl Bk 02444	617 735-8787
	Carter J M 3410 Columbia Rd S 02136	617 464-1040
	Carter J M Ornamental Ironworks 40 Franklin Waterbury 02172	617 324-5353
	Carter J Neal Co 40 Newbury St 02118	617 442-1775
	Carter James 1573 Cambridge St Cam 02138	617 492-1214
	James 422 Foster Av Roxbury 02132	617 739-2193
	James 31 East Star Rd Cambridge 02141	617 876-8841
	Jane L 34 Rosbury Rd Mt 02136	617 361-0773
	Janine 14 Adams Rd Newton 02459	617 564-0435
	Jenny 41 Warren Av Sun 02145	617 426-5994
	John 11 Mansfield Mt 02134	617 987-2163
	John 207 Summer St 02125	617 423-4334
	John 40 Westford Der 02125	617 282-1235
	Juanes O 129 A Summit Av 02131	617 734-6109
	K 29 Fenwick Av East Boston 02128	617 565-4856
	K 7 Fenwick Der 02125	617 282-1593
	Carter Nellie E 323 Marchant Av Mt 02115	617 267-6483
	Nicholas S F 115 Randolph Av Mt 02136	617 698-5307
	Nick 21 Furlow Bay 02114	617 267-5222
	Nick & Debbie 136 Hermit Rd Newton 02459	617 527-0480
	Norman G 38 Chickadee Rd Der 02125	617 822-1201
	P 46 Cranston Pl 02115	617 437-4754
	P E 501 E South St 02037	617 268-4213
	P L 44 Hutchings Bay 02115	617 427-9170
	P R 91 Brewer Jan 02138	617 968-8692
	Paul & Constance 114 Beacon Av W Mt 02110	617 325-3034
	Paul E 501 E South St 02037	617 268-4546
	Paul M 27 Union St 02135	617 787-2115
	Carter Pike Driving Inc 27 Beaver Ct Framingham 02702	Wellesley Tpk-781.235-8488
	Carter Prudence 40 Franklin Waterbury 02172	617 393-3782
	Prudence 40 Franklin Waterbury 02172	617 926-7063
	Reginald 100 Fenwick Cambridge 02124	617 541-2843
	Renee & Andrew 10 Walnut St 02118	617 720-3765
	Carter Rice Donald Building Services Publishing 163 Main Wilmington 01887 Toll Free-800 J & Thom.....800 638-1671	
	Carl E. Anderson & Fred 613 Main Wilmington Toll Free-800 J & Thom.....800 619-7447	
	Toll Free-800 J & Thom.....800 648-7447	
	Reynolds 413 Main Wilmington 01887978 988-7447	
	Ingalls Circle 163 Main Wilmington 01887800 638-1673	
	Carter Richard 2075 Cambridge Av Brighton 02215	617 982-0836
	Richard A 47 West Vernon St 02136	617 566-7293
	Carter Richard A MD 120 Fenwick Av Sun 02145	617 267-0710
	Carter Richard K 123 Mercer St 02127	617 268-0468
	Robert L 175 Rockwood Av Cam 02141	617 864-1535
	Roger 130 Stoughton St 02114	617 424-6148
	Royce & Andrew 18 Salisbury Cha 02129	617 491-6115
	Royce 18 Salisbury Cha 02129	617 241-9418



A New Strategy

Make it easy to choose best clustering from millions of choices

A New Strategy

Make it easy to choose best clustering from millions of choices

- 1 **Code text as numbers** (in one *or more* of several ways)

A New Strategy

Make it easy to choose best clustering from millions of choices

- 1 **Code text as numbers** (in one *or more* of several ways)
- 2 **Apply all clustering methods we can find** to the data — each representing different (unstated) substantive assumptions (<15 mins)

A New Strategy

Make it easy to choose best clustering from millions of choices

- 1 **Code text as numbers** (in one *or more* of several ways)
- 2 **Apply all clustering methods we can find** to the data — each representing different (unstated) substantive assumptions (<15 mins)
- 3 (Too much for a person to understand, but organization will help)

A New Strategy

Make it easy to choose best clustering from millions of choices

- 1 Code text as numbers (in one or more of several ways)
- 2 Apply all clustering methods we can find to the data — each representing different (unstated) substantive assumptions (<15 mins)
- 3 (Too much for a person to understand, but organization will help)
- 4 Develop an application-independent distance metric between clusterings, a metric space of clusterings, and a 2-D projection

A New Strategy

Make it easy to choose best clustering from millions of choices

- 1 **Code text as numbers** (in one *or more* of several ways)
- 2 **Apply all clustering methods we can find** to the data — each representing different (unstated) substantive assumptions (<15 mins)
- 3 (Too much for a person to understand, but organization will help)
- 4 Develop an **application-independent distance metric** between clusterings, a **metric space of clusterings**, and a **2-D projection**
- 5 “**Local cluster ensemble**” creates a new clustering at any point, based on weighted average of nearby clusterings

A New Strategy

Make it easy to choose best clustering from millions of choices

- 1 **Code text as numbers** (in one *or more* of several ways)
- 2 **Apply all clustering methods we can find** to the data — each representing different (unstated) substantive assumptions (<15 mins)
- 3 (Too much for a person to understand, but organization will help)
- 4 Develop an **application-independent distance metric** between clusterings, a **metric space of clusterings**, and a **2-D projection**
- 5 “**Local cluster ensemble**” creates a new clustering at any point, based on weighted average of nearby clusterings
- 6 A new **animated visualization** to explore the space of clusterings (smoothly morphing from one into others)

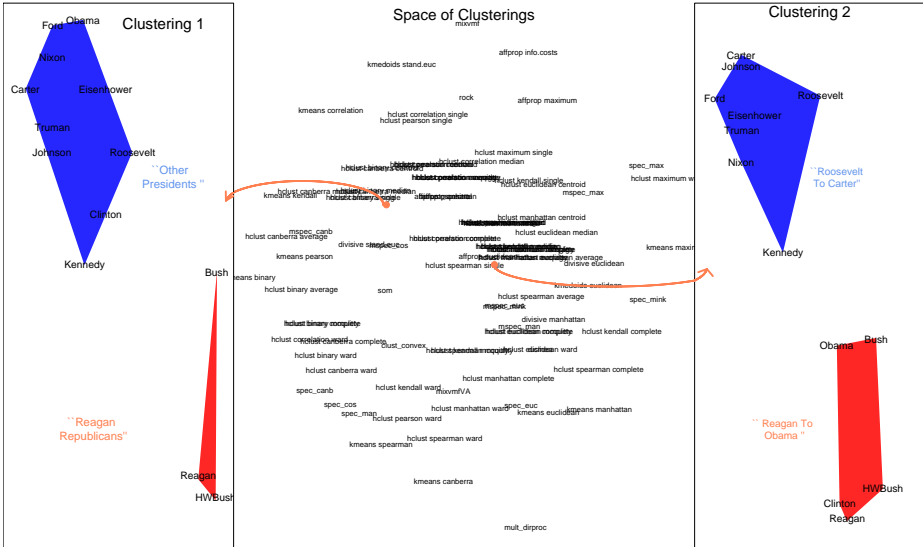
A New Strategy

Make it easy to choose best clustering from millions of choices

- 1 **Code text as numbers** (in one or more of several ways)
- 2 **Apply all clustering methods we can find** to the data — each representing different (unstated) substantive assumptions (<15 mins)
- 3 (Too much for a person to understand, but organization will help)
- 4 Develop an **application-independent distance metric** between clusterings, a **metric space of clusterings**, and a **2-D projection**
- 5 “**Local cluster ensemble**” creates a new clustering at any point, based on weighted average of nearby clusterings
- 6 A new **animated visualization** to explore the space of clusterings (smoothly morphing from one into others)
- 7 **↔ Millions of clusterings, easily comprehended** (takes about 10-15 minutes to choose a clustering with insight)

Many Thousands of Clusterings, Sorted & Organized

You choose one (or more), based on insight, discovery, useful information, . . .



Application-Independent Distance Metric: Axioms

Application-Independent Distance Metric: Axioms

- Metric based on 3 assumptions

Application-Independent Distance Metric: Axioms

- Metric based on 3 assumptions
 - ① Distance between clusterings: a function of the **pairwise document agreements** (pairwise agreements \Rightarrow triples, quadruples, etc.)

Application-Independent Distance Metric: Axioms

- Metric based on 3 assumptions
 - ① Distance between clusterings: a function of the **pairwise document agreements** (pairwise agreements \Rightarrow triples, quadruples, etc.)
 - ② **Invariance**: Distance is invariant to the number of documents (for any fixed number of clusters)

Application-Independent Distance Metric: Axioms

- Metric based on 3 assumptions
 - ① Distance between clusterings: a function of the **pairwise document agreements** (pairwise agreements \Rightarrow triples, quadruples, etc.)
 - ② **Invariance**: Distance is invariant to the number of documents (for any fixed number of clusters)
 - ③ **Scale**: the maximum distance is set to $\log(\text{num clusters})$

Application-Independent Distance Metric: Axioms

- Metric based on 3 assumptions
 - ① Distance between clusterings: a function of the **pairwise document agreements** (pairwise agreements \Rightarrow triples, quadruples, etc.)
 - ② **Invariance**: Distance is invariant to the number of documents (for any fixed number of clusters)
 - ③ **Scale**: the maximum distance is set to $\log(\text{num clusters})$
- \rightsquigarrow **Only one measure satisfies all three** (the “variation of information”)

Application-Independent Distance Metric: Axioms

- Metric based on 3 assumptions
 - ① Distance between clusterings: a function of the **pairwise document agreements** (pairwise agreements \Rightarrow triples, quadruples, etc.)
 - ② **Invariance**: Distance is invariant to the number of documents (for any fixed number of clusters)
 - ③ **Scale**: the maximum distance is set to $\log(\text{num clusters})$
- \rightsquigarrow **Only one measure satisfies all three** (the “variation of information”)
- Meila (2007): derives same metric using different axioms (lattice theory)

Evaluating the Performance of Our Method

Evaluating the Performance of Our Method

- Goals:

Evaluating the Performance of Our Method

- Goals:
 - **Validate Claim:** computer-assisted conceptualization outperforms human conceptualization

Evaluating the Performance of Our Method

- Goals:
 - **Validate Claim:** computer-assisted conceptualization outperforms human conceptualization
 - **Demonstrate:** new experimental designs for cluster evaluation

Evaluating the Performance of Our Method

- Goals:
 - **Validate Claim**: computer-assisted conceptualization outperforms human conceptualization
 - **Demonstrate**: new experimental designs for cluster evaluation
 - **Inject human judgement**: relying on insights from survey research

Evaluating the Performance of Our Method

- Goals:
 - **Validate Claim**: computer-assisted conceptualization outperforms human conceptualization
 - **Demonstrate**: new experimental designs for cluster evaluation
 - **Inject human judgement**: relying on insights from survey research
- We now present three evaluations

Evaluating the Performance of Our Method

- Goals:
 - **Validate Claim**: computer-assisted conceptualization outperforms human conceptualization
 - **Demonstrate**: new experimental designs for cluster evaluation
 - **Inject human judgement**: relying on insights from survey research
- We now present three evaluations
 - Quality \Rightarrow RA coders

Evaluating the Performance of Our Method

- Goals:
 - **Validate Claim**: computer-assisted conceptualization outperforms human conceptualization
 - **Demonstrate**: new experimental designs for cluster evaluation
 - **Inject human judgement**: relying on insights from survey research
- We now present three evaluations
 - Quality \Rightarrow RA coders
 - Informative discoveries \Rightarrow Experienced scholars analyzing texts

Evaluating the Performance of Our Method

- Goals:
 - **Validate Claim**: computer-assisted conceptualization outperforms human conceptualization
 - **Demonstrate**: new experimental designs for cluster evaluation
 - **Inject human judgement**: relying on insights from survey research
- We now present three evaluations
 - Quality \Rightarrow RA coders
 - Informative discoveries \Rightarrow Experienced scholars analyzing texts
 - Discovery \Rightarrow You're the judge

Evaluation 1: Cluster Quality

Evaluation 1: Cluster Quality

- What Are Humans Good For?

Evaluation 1: Cluster Quality

- **What Are Humans Good For?**
 - They can't: keep many documents & clusters in their head

Evaluation 1: Cluster Quality

- **What Are Humans Good For?**
 - They can't: keep many documents & clusters in their head
 - They can: compare two documents at a time

Evaluation 1: Cluster Quality

- **What Are Humans Good For?**
 - They can't: keep many documents & clusters in their head
 - They can: compare two documents at a time
 - \implies Cluster quality evaluation: human judgement of document pairs

Evaluation 1: Cluster Quality

- **What Are Humans Good For?**
 - They can't: keep many documents & clusters in their head
 - They can: compare two documents at a time
 - \implies Cluster quality evaluation: human judgement of document pairs
- **Experimental Design to Assess Cluster Quality**

Evaluation 1: Cluster Quality

- **What Are Humans Good For?**
 - They can't: keep many documents & clusters in their head
 - They can: compare two documents at a time
 - \implies Cluster quality evaluation: human judgement of document pairs
- **Experimental Design to Assess Cluster Quality**
 - automated visualization to choose one clustering

Evaluation 1: Cluster Quality

- **What Are Humans Good For?**
 - They can't: keep many documents & clusters in their head
 - They can: compare two documents at a time
 - \implies Cluster quality evaluation: human judgement of document pairs
- **Experimental Design to Assess Cluster Quality**
 - automated visualization to choose one clustering
 - many pairs of documents

Evaluation 1: Cluster Quality

- **What Are Humans Good For?**
 - They can't: keep many documents & clusters in their head
 - They can: compare two documents at a time
 - \implies Cluster quality evaluation: human judgement of document pairs
- **Experimental Design to Assess Cluster Quality**
 - automated visualization to choose one clustering
 - many pairs of documents
 - for coders: (1) unrelated, (2) loosely related, (3) closely related

Evaluation 1: Cluster Quality

- **What Are Humans Good For?**
 - They can't: keep many documents & clusters in their head
 - They can: compare two documents at a time
 - \implies Cluster quality evaluation: human judgement of document pairs
- **Experimental Design to Assess Cluster Quality**
 - automated visualization to choose one clustering
 - many pairs of documents
 - for coders: (1) unrelated, (2) loosely related, (3) closely related
 - Quality = mean(within cluster) - mean(between clusters)

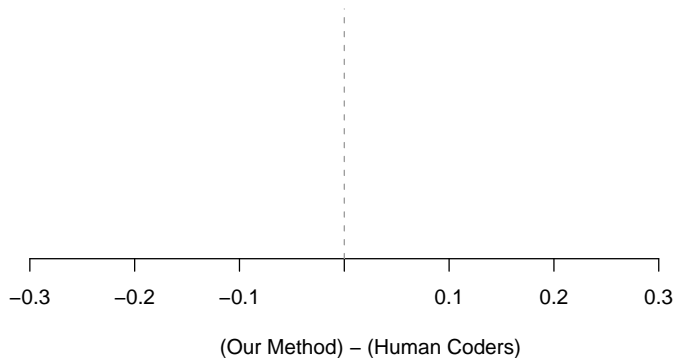
Evaluation 1: Cluster Quality

- **What Are Humans Good For?**
 - They can't: keep many documents & clusters in their head
 - They can: compare two documents at a time
 - \implies Cluster quality evaluation: human judgement of document pairs
- **Experimental Design to Assess Cluster Quality**
 - automated visualization to choose one clustering
 - many pairs of documents
 - for coders: (1) unrelated, (2) loosely related, (3) closely related
 - Quality = mean(within cluster) - mean(between clusters)
 - **Bias results against ourselves by not letting evaluators choose clustering**

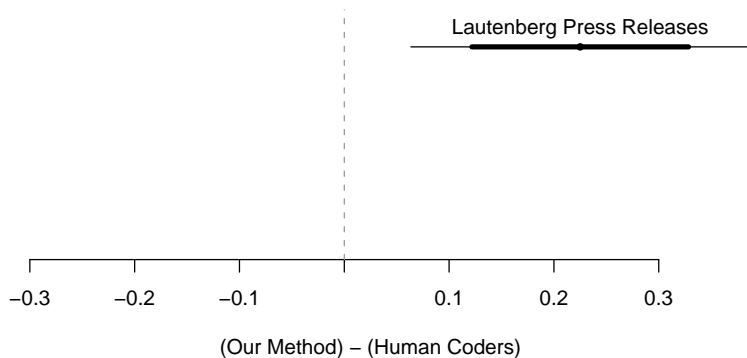
Evaluation 1: Cluster Quality

- **What Are Humans Good For?**
 - They can't: keep many documents & clusters in their head
 - They can: compare two documents at a time
 - \implies Cluster quality evaluation: human judgement of document pairs
- **Experimental Design to Assess Cluster Quality**
 - automated visualization to choose one clustering
 - many pairs of documents
 - for coders: (1) unrelated, (2) loosely related, (3) closely related
 - Quality = mean(within cluster) - mean(between clusters)
 - **Bias results against ourselves by not letting evaluators choose clustering**

Evaluation 1: Cluster Quality

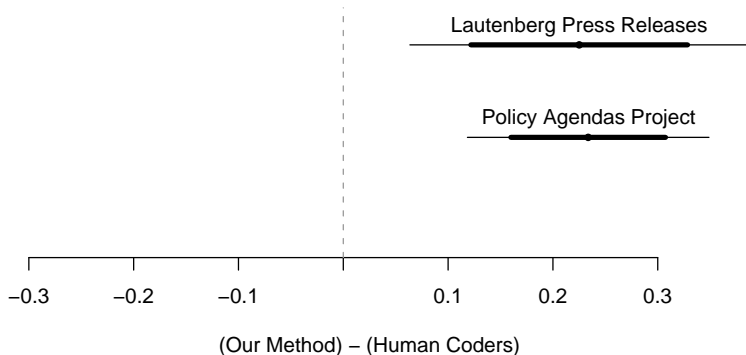


Evaluation 1: Cluster Quality



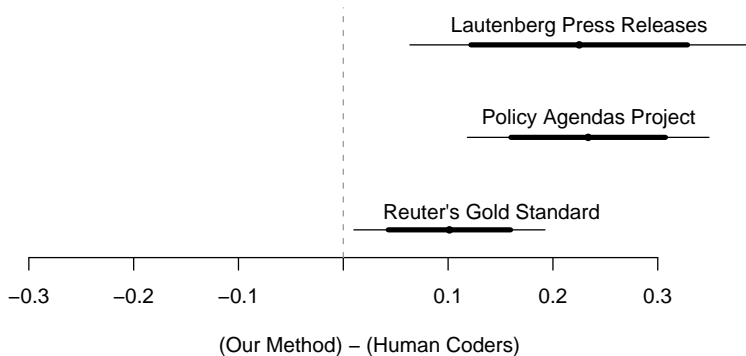
Lautenberg: 200 Senate Press Releases (appropriations, economy, education, tax, veterans, ...)

Evaluation 1: Cluster Quality



Policy Agendas: 213 quasi-sentences from Bush's State of the Union (agriculture, banking & commerce, civil rights/liberties, defense, ...)

Evaluation 1: Cluster Quality



Reuter's: financial news (trade, earnings, copper, gold, coffee, . . .); "gold standard" for supervised learning studies

Evaluation 2: More Informative Discoveries

Evaluation 2: More Informative Discoveries

- Found 2 scholars analyzing lots of textual data for their work

Evaluation 2: More Informative Discoveries

- Found 2 scholars analyzing lots of textual data for their work
- Created 6 clusterings:

Evaluation 2: More Informative Discoveries

- Found 2 scholars analyzing lots of textual data for their work
- Created 6 clusterings:
 - 2 clusterings selected with our method (**biased** against us)

Evaluation 2: More Informative Discoveries

- Found 2 scholars analyzing lots of textual data for their work
- Created 6 clusterings:
 - 2 clusterings selected with our method (**biased** against us)
 - 2 clusterings from each of 2 other methods (varying tuning parameters)

Evaluation 2: More Informative Discoveries

- Found 2 scholars analyzing lots of textual data for their work
- Created 6 clusterings:
 - 2 clusterings selected with our method (**biased** against us)
 - 2 clusterings from each of 2 other methods (varying tuning parameters)
- Created info packet on each clustering (for each cluster: exemplar document, automated content summary)

Evaluation 2: More Informative Discoveries

- Found 2 scholars analyzing lots of textual data for their work
- Created 6 clusterings:
 - 2 clusterings selected with our method (**biased** against us)
 - 2 clusterings from each of 2 other methods (varying tuning parameters)
- Created info packet on each clustering (for each cluster: exemplar document, automated content summary)
- Asked for $\binom{6}{2}=15$ pairwise comparisons

Evaluation 2: More Informative Discoveries

- Found 2 scholars analyzing lots of textual data for their work
- Created 6 clusterings:
 - 2 clusterings selected with our method (**biased** against us)
 - 2 clusterings from each of 2 other methods (varying tuning parameters)
- Created info packet on each clustering (for each cluster: exemplar document, automated content summary)
- Asked for $\binom{6}{2}=15$ pairwise comparisons
- User chooses \Rightarrow only care about the one clustering that wins

Evaluation 2: More Informative Discoveries

- Found 2 scholars analyzing lots of textual data for their work
- Created 6 clusterings:
 - 2 clusterings selected with our method (**biased** against us)
 - 2 clusterings from each of 2 other methods (varying tuning parameters)
- Created info packet on each clustering (for each cluster: exemplar document, automated content summary)
- Asked for $\binom{6}{2}=15$ pairwise comparisons
- User chooses \Rightarrow only care about the one clustering that wins
- Both cases a Condorcet winner:

Evaluation 2: More Informative Discoveries

- Found 2 scholars analyzing lots of textual data for their work
- Created 6 clusterings:
 - 2 clusterings selected with our method (**biased** against us)
 - 2 clusterings from each of 2 other methods (varying tuning parameters)
- Created info packet on each clustering (for each cluster: exemplar document, automated content summary)
- Asked for $\binom{6}{2}=15$ pairwise comparisons
- User chooses \Rightarrow only care about the one clustering that wins
- Both cases a Condorcet winner:

“Immigration”:

Our Method 1 \rightarrow vMF 1 \rightarrow vMF 2 \rightarrow Our Method 2 \rightarrow K-Means 1 \rightarrow K-Means 2

Evaluation 2: More Informative Discoveries

- Found 2 scholars analyzing lots of textual data for their work
- Created 6 clusterings:
 - 2 clusterings selected with our method (**biased** against us)
 - 2 clusterings from each of 2 other methods (varying tuning parameters)
- Created info packet on each clustering (for each cluster: exemplar document, automated content summary)
- Asked for $\binom{6}{2}=15$ pairwise comparisons
- User chooses \Rightarrow only care about the one clustering that wins
- Both cases a Condorcet winner:

“Immigration”:

Our Method 1 \rightarrow vMF 1 \rightarrow vMF 2 \rightarrow Our Method 2 \rightarrow K-Means 1 \rightarrow K-Means 2

“Genetic testing”:

Our Method 1 \rightarrow {Our Method 2, K-Means 1, K-means 2} \rightarrow Dir Proc. 1 \rightarrow Dir Proc. 2

Evaluation 3: What Do Members of Congress Do?

Evaluation 3: What Do Members of Congress Do?

- David Mayhew's (1974) famous typology

Evaluation 3: What Do Members of Congress Do?

- David Mayhew's (1974) famous typology
 - Advertising

Evaluation 3: What Do Members of Congress Do?

- David Mayhew's (1974) famous typology
 - Advertising
 - Credit Claiming

Evaluation 3: What Do Members of Congress Do?

- David Mayhew's (1974) famous typology
 - Advertising
 - Credit Claiming
 - Position Taking

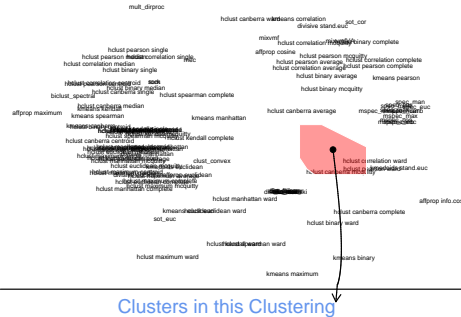
Evaluation 3: What Do Members of Congress Do?

- David Mayhew's (1974) famous typology
 - Advertising
 - Credit Claiming
 - Position Taking
- Data: 200 press releases from Frank Lautenberg's office (D-NJ)

Evaluation 3: What Do Members of Congress Do?

- David Mayhew's (1974) famous typology
 - Advertising
 - Credit Claiming
 - Position Taking
- Data: 200 press releases from Frank Lautenberg's office (D-NJ)
- Apply our method

Example Discovery



Credit Claiming, Legislation:
“As the Senate begins its recess, Senator Frank Lautenberg today pointed to a string of victories in Congress on his legislative agenda during this work period”

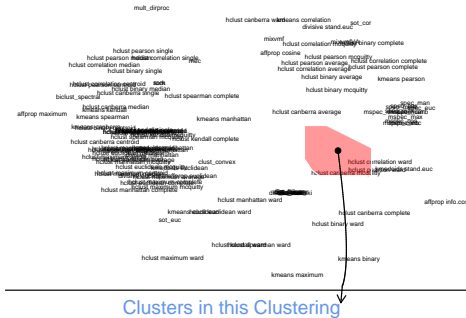


Credit Claiming
Pork



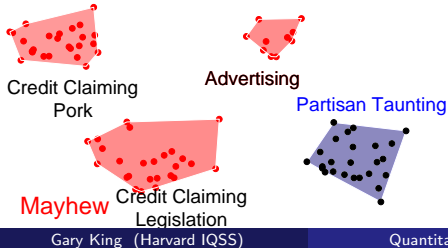
Mayhew
Credit Claiming
Legislation
Gary King (Harvard IQSS)

Example Discovery: Partisan Taunting



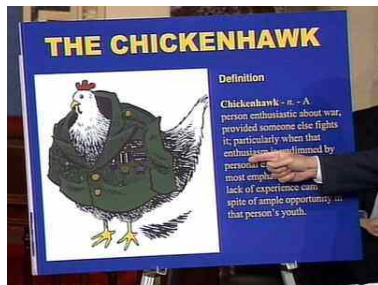
Definition: Explicit, public, and negative attacks on another political party or its members

Taunting ruins deliberation



In Sample Illustration of Partisan Taunting

Taunting ruins deliberation

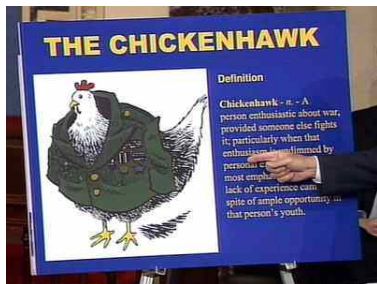


Sen. Lautenberg
on Senate Floor
4/29/04

- "Senator Lautenberg Blasts Republicans as 'Chicken Hawks' " [Government Oversight]

In Sample Illustration of Partisan Taunting

Taunting ruins deliberation



Sen. Lautenberg
on Senate Floor
4/29/04

- "Senator Lautenberg Blasts Republicans as 'Chicken Hawks' " [Government Oversight]
- "The scopes trial took place in 1925. Sadly, President Bush's veto today shows that we haven't progressed much since then" [Healthcare]

Taunting ruins deliberation



Sen. Lautenberg
on Senate Floor
4/29/04

- "Senator Lautenberg Blasts Republicans as 'Chicken Hawks' " [Government Oversight]
- "The scopes trial took place in 1925. Sadly, President Bush's veto today shows that we haven't progressed much since then" [Healthcare]
- "Every day the House Republicans dragged this out was a day that made our communities less safe." [Homeland Security]

Out of Sample Confirmation of Partisan Taunting

- Discovered using 200 press releases; 1 senator.

Out of Sample Confirmation of Partisan Taunting

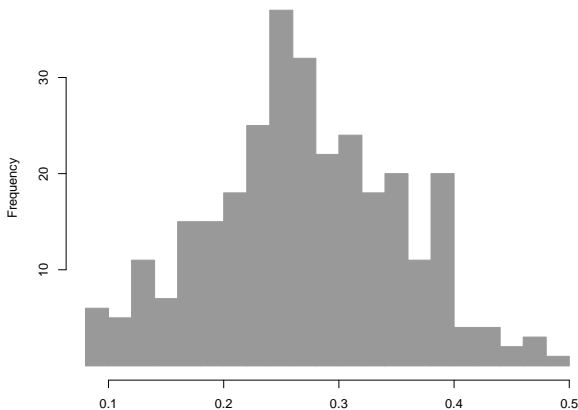
- Discovered using 200 press releases; 1 senator.
- Confirmed using 64,033 press releases; 301 senator-years.

Out of Sample Confirmation of Partisan Taunting

- Discovered using 200 press releases; 1 senator.
- Confirmed using 64,033 press releases; 301 senator-years.
- Apply supervised learning method: measure **proportion of press releases** a senator taunts other party

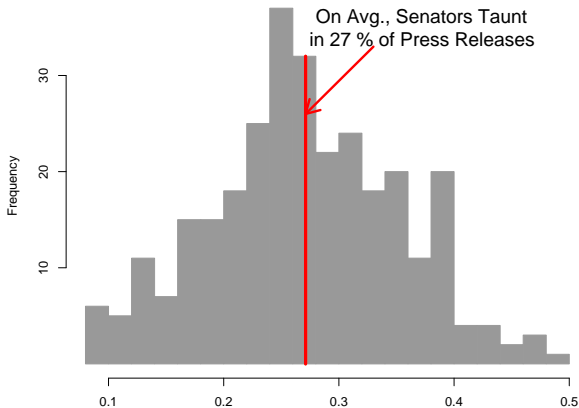
Out of Sample Confirmation of Partisan Taunting

- Discovered using 200 press releases; 1 senator.
- Confirmed using 64,033 press releases; 301 senator-years.
- Apply supervised learning method: measure **proportion of press releases** a senator taunts other party

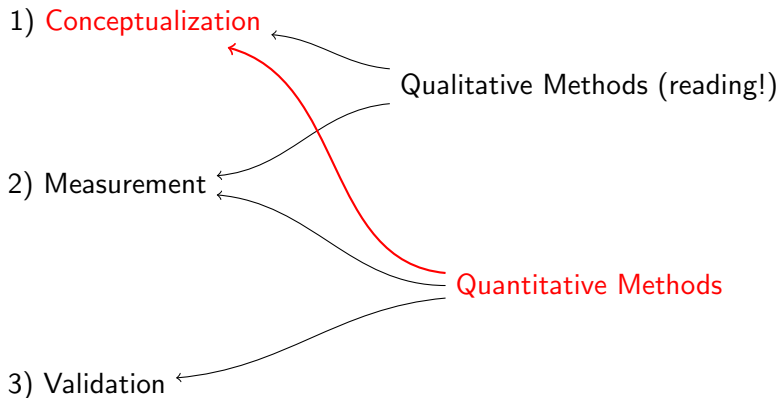


Out of Sample Confirmation of Partisan Taunting

- Discovered using 200 press releases; 1 senator.
- Confirmed using 64,033 press releases; 301 senator-years.
- Apply supervised learning method: measure **proportion of press releases** a senator taunts other party

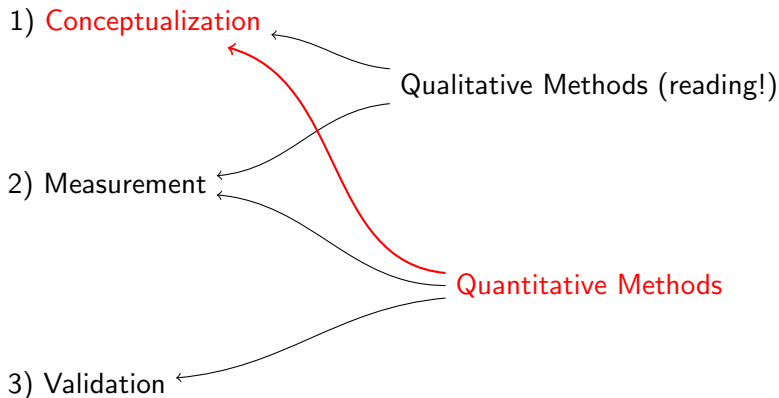


Advancing the Objective of Discovery



Quantitative methods for conceptualization: aiding **discovery**

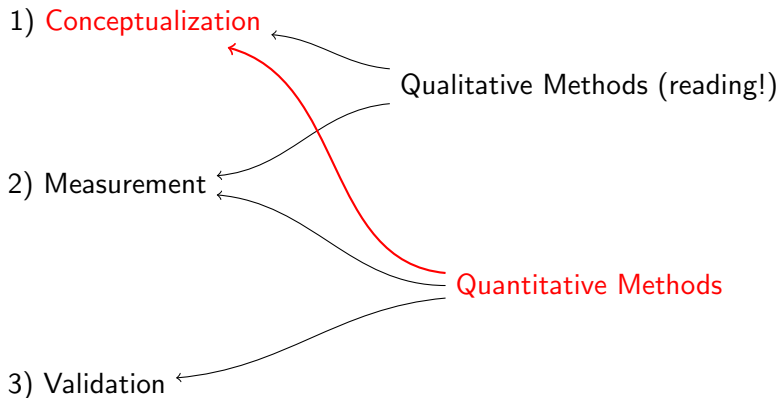
Advancing the Objective of Discovery



Quantitative methods for conceptualization: aiding **discovery**

- Few formal methods designed explicitly for conceptualization

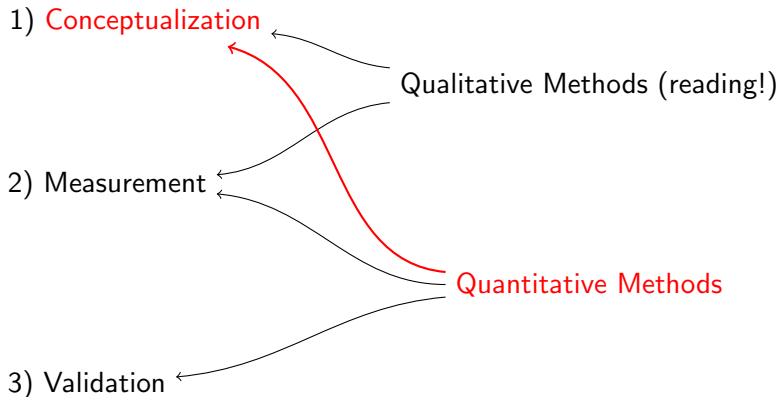
Advancing the Objective of Discovery



Quantitative methods for conceptualization: aiding **discovery**

- Few formal methods designed explicitly for conceptualization
- **Belittled**: “Tom Swift and His Electric Factor Analysis Machine” (Armstrong 1967)

Advancing the Objective of Discovery



Quantitative methods for conceptualization: aiding **discovery**

- Few formal methods designed explicitly for conceptualization
- **Belittled**: “Tom Swift and His Electric Factor Analysis Machine” (Armstrong 1967)
- Evaluation methods measure progress in discovery

For more information (on adding zooming out to the human ability to zoom in)

<http://GKing.Harvard.edu>