

A General Purpose Computer-Assisted Clustering Methodology

Gary King

Institute for Quantitative Social Science
Harvard University

Talk at University of Massachusetts, Amherst, 10/28/2010

Joint work with Justin Grimmer (Harvard \rightsquigarrow Stanford)

A Method for Computer Assisted Conceptualization

- 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 Computer Assisted Conceptualization

- 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 Computer Assisted Conceptualization

- 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
- Main goal: Switch from **Fully Automated** to **Computer Assisted**

A Method for Computer Assisted Conceptualization

- 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
- Main goal: Switch from **Fully Automated** to **Computer Assisted**
- (We focus on clustering texts; methods apply more broadly)

What's Hard about Clustering?

What's Hard about Clustering?

(aka Why Johnny Can't Classify)

What's Hard about Clustering?

(aka Why Johnny Can't Classify)

- Clustering seems easy; its not!

What's Hard about Clustering?

(aka Why Johnny Can't Classify)

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

What's Hard about Clustering?

(aka Why Johnny Can't Classify)

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

What's Hard about Clustering?

(aka Why Johnny Can't Classify)

- 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)

What's Hard about Clustering?

(aka Why Johnny Can't Classify)

- 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)
- $Bell(5) = 52$

What's Hard about Clustering?

(aka Why Johnny Can't Classify)

- 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)
- $Bell(5) = 52$
- $Bell(100) \approx$

What's Hard about Clustering?

(aka Why Johnny Can't Classify)

- 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

What's Hard about Clustering?

(aka Why Johnny Can't Classify)

- 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!

What's Hard about Clustering?

(aka Why Johnny Can't Classify)

- 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!
- Fully automated algorithms can help, but which ones?

The Problem with Fully Automated Clustering

The Problem with Fully Automated Clustering

- **The (Impossible) Goal:** optimal, fully automated, application-independent cluster analysis

The Problem with Fully Automated Clustering

- **The (Impossible) Goal:** optimal, fully automated, application-independent cluster analysis
- **No free lunch theorem:** every possible clustering method performs equally well on average over all possible substantive applications

The Problem with Fully Automated Clustering

- **The (Impossible) Goal:** optimal, fully automated, application-independent cluster analysis
- **No free lunch theorem:** every possible clustering method performs equally well on average over all possible substantive applications
- Existing methods:

The Problem with Fully Automated Clustering

- **The (Impossible) Goal:** optimal, fully automated, application-independent cluster analysis
- **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,...

The Problem with Fully Automated Clustering

- **The (Impossible) Goal:** optimal, fully automated, application-independent cluster analysis
- **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

The Problem with Fully Automated Clustering

- **The (Impossible) Goal:** optimal, fully automated, application-independent cluster analysis
- **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 Problem with Fully Automated Clustering

- **The (Impossible) Goal:** optimal, fully automated, application-independent cluster analysis
- **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**

The Problem with Fully Automated Clustering

- **The (Impossible) Goal:** optimal, fully automated, application-independent cluster analysis
- **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

The Problem with Fully Automated Clustering

- **The (Impossible) Goal:** optimal, fully automated, application-independent cluster analysis
- **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:** full automation requires more information

The Problem with Fully Automated Clustering

- **The (Impossible) Goal:** optimal, fully automated, application-independent cluster analysis
- **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:** full automation requires more information
- No surprise: everyone's tried cluster analysis; very few are satisfied

Switch from Fully Automated to Computer Assisted

Switch from Fully Automated to Computer Assisted

- **Fully Automated Clustering** may succeed sometimes, but fails in general: too hard to understand when each model applies

Switch from Fully Automated to Computer Assisted

- **Fully Automated Clustering** may succeed sometimes, but fails in general: too hard to understand when each model applies
- An alternative: **Computer-Assisted Clustering**

Switch from Fully Automated to Computer Assisted

- **Fully Automated Clustering** may succeed sometimes, but fails in general: too hard to understand when each model applies
- An alternative: **Computer-Assisted Clustering**
 - Create long list of clusterings; choose the best

Switch from Fully Automated to Computer Assisted

- **Fully Automated Clustering** may succeed sometimes, but fails in general: too hard to understand when each model applies
- An alternative: **Computer-Assisted Clustering**
 - Create long list of clusterings; choose the best
 - Too hard for mere humans!

Switch from Fully Automated to Computer Assisted

- **Fully Automated Clustering** may succeed sometimes, but fails in general: too hard to understand when each model applies
- An alternative: **Computer-Assisted Clustering**
 - Create long list of clusterings; choose the best
 - Too hard for mere humans!
 - An **organized** list will make the search possible

Switch from Fully Automated to Computer Assisted

- **Fully Automated Clustering** may succeed sometimes, but fails in general: too hard to understand when each model applies
- An alternative: **Computer-Assisted Clustering**
 - Create long list of clusterings; choose the best
 - Too hard for mere humans!
 - An **organized** list will make the search possible
 - Insight: Many clusterings are perceptually identical

Switch from Fully Automated to Computer Assisted

- **Fully Automated Clustering** may succeed sometimes, but fails in general: too hard to understand when each model applies
- An alternative: **Computer-Assisted Clustering**
 - Create long list of clusterings; choose the best
 - Too hard for mere humans!
 - An **organized** list will make the search possible
 - Insight: Many clusterings are perceptually identical
 - E.g.,: consider two clusterings that differ only because one document (of 10,000) moves from category 5 to 6

Switch from Fully Automated to Computer Assisted

- **Fully Automated Clustering** may succeed sometimes, but fails in general: too hard to understand when each model applies
- An alternative: **Computer-Assisted Clustering**
 - Create long list of clusterings; choose the best
 - Too hard for mere humans!
 - An **organized** list will make the search possible
 - Insight: Many clusterings are perceptually identical
 - E.g.,: consider two clusterings that differ only because one document (of 10,000) moves from category 5 to 6
- **The Question: How to organize all those clusterings?**

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	Carter F 34 Hibiscus Bldg 02133..... 617 327-1105	Carter Nella E 323 Main St 02115..... 617 267-6483	
Cartagena Lydia 28 Sweet Box 02131..... 617 323-7639	Faye & Ricky 207 Columbia Ave Box 02136..... 617 437-7331	Nicholas S F 115 Randolph Ave Mill 02186..... 617 698-5307	
Cartagena Avish F Pleasant Hill 02139..... 617 442-9780	Francis S 134 Yankov W Ave 02132..... 617 323-6781	Nick 21 Farwell Box 02114..... 617 267-5222	
B Had 02134..... 617 361-5253	Franklin & Anne 705 Mt Auburn Cam 02138..... 617 354-0798	Nick & Debbi 196 Herold Rd Newton 02459..... 617 527-0480	
Jessica 50 Decatur Cha 02129..... 617 241-0152	Fred 41 Howland Elm 02136..... 617 524-3078	Nicole..... 617 698-0713	
Luzmila 124 Harvard Cam 02136..... 617 491-5621	Fred 16 Howland Ave Mill 02136..... 617 698-1343	Norman G 38 Chickareed Dr 02125..... 617 822-1201	
M 95 Howe Box 02132..... 617 323-9713	G & B 8 Vardon Dr 02134..... 617 434-8966	P 40 Cranford Pl Box 02135..... 617 437-4754	
Melvin 503 Green Cam 02129..... 617 576-1061	G T 27 Franklin Ave Sun 02145..... 617 623-7121	P E 501 E South S Box 02137..... 617 268-8213	
Carte Nicholas 18 Appleton Boston 02114..... 617 695-6996	George 125 Madison Box 02134..... 617 367-9548	P E 14 Hutchings Box 02131..... 617 427-9170	
Carlencio 0 4 Bedford Box 02133..... 617 338-9219	Carter Hillside Assoc/Am 107 S Street Box 02111..... 617 456-1689	Paul & Constance 114 Amesbury W Box 02131..... 617 325-2036	
Carten Thos Jr Sr & Claire 1 Franklin Hill Mt 02136..... 617 698-6163	Carter Harry F 30 Bayview Rd Rt W Box 02132..... 617 325-5465	Paul M 501 E South S S Box 02137..... 617 268-4546	
17 445-5116	Carte Hide Co Inc 26 Irving St 02133..... 617 542-7987	Paul M 27 Crown Bk 02139..... 617 787-2115	
17 822-2992	A Heber 617 442-5230	Carter Pile Driving Inc 27 Amesbury Ct Frankington 02102..... Woblesley Tpk 781.235-0488	
17 427-5712	Carter Hilary 41 Harvey Cam 02148..... 617 876-2750	Carter Prudence 40 Franklin Waterman 02127..... 617 393-3782	
17 569-2698	Horace 301 Walnut St Roxbury 02119..... 617 442-5307	Prudence 40 Franklin Waterman 02127..... 617 926-7063	
17 667-5190	Howard Jr 28 Nona Drive Box 02118..... 617 445-5532	Roginald 106 Brookview Dorchester 02122..... 617 541-2843	
17 569-1412	J Dan..... 617 354-2658	Renee & Andrew 106 Brookview Dorchester 02122..... 617 720-3765	
17 338-9110	J 31 Chatham Box 02144..... 617 232-7990		
17 825-1993	J 538 Harvard Box 02144..... 617 730-9483		
17 296-1593	J 775 The Pines West Roxbury 02132..... 617 323-5374		
17 670-2078	Jacques J Jacques MD 1 Brookline Pl Box 02144..... 617 735-8787		
17 621-9001	Carter J M 3410 Columbia Rd S Box 02137..... 617 464-1040		
17 296-4725	Carter J M Ornamental Ironworks 200 Walnut St Roxbury 02119..... 617 442-5307		
17 542-1521	Carter J Neal Co 40 Howland Elm 02136..... 617 442-1775		
17 364-5232	Carter James 157 Cambridge St Cam 02138..... 617 492-1214		
17 541-5649	Carter J Broadcasting Co 30 Park Pl Box 02134..... 617 423-0210		
17 739-2662	Carter & Bussac Consultants Inc 73 East St Cam 02141..... 617 225-0200		
17 879-0030	Carter C 200 Commonwealth Ave 02135..... 617 782-2118		
17 436-1511	C 218 Harvard Ave East Boston 02128..... 617 569-1545		
17 569-4119	C 109 Harvard Cam 02136..... 617 491-4822		
800 569-4782	C & M 43 Bernham Jan 02136..... 617 524-9558		

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
17 447-4101	Cartagema Lydia 28 Sweet Briar Dr 02131	617 323-7639
100 257-9961	Cartagema Avish F Beach Rd 02139	617 442-9780
	B Had 02134	617 361-5253
17 566-1282	Cartage J 50 Decatur Cha 02129	617 241-0152
17 364-5188	Cartage L 124 Harvard Cam 02138	617 491-5621
	M 95 Howe St 02136	617 323-9713
361-0380	Carte Nicholas Milton 503 Green Cam 02139	617 576-1061
17 566-4548	Carte Nicholas 18 Appleton Boston 02114	617 695-6996
	Cartagema O 4 Harvard Boy 02138	617 338-9219
17 628-8248	Carten Thos J Sr & Claire 1 Furlong St Mt 02136	617 698-6163
17 445-5116	Carton & Kullback 50 Thompson Ln Mt 02136	617 696-6919
17 822-2962	Cartor A Inc 02133	617 229-2257
17 427-5712	A Nader A 22 Bethune 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
17 569-1417	Adams 301 Carter St Mt 02136	617 698-9074
17 338-9110	Adams P 40 Market Cambridge 02139	617 945-2711
17 825-9195	Adams P 42 Mt St 02138	617 625-7623
17 296-1593	Adams P 1161 Beacon St 02144	617 739-1022
17 670-2078	B E 10 Gladstone Av Mt 02136	617 536-6229
17 621-9001	Adams P Turk New England Medical Center Inc 02111	617 296-6911
17 296-4725	Cartor Beckey Inc 02134	617 636-0951
17 542-1521	Bernard J 371 Newbury Boston 02116	617 523-4368
17 364-5232	Bibbith 25 Midway Dr 02134	617 567-9430
17 541-5649	Bibbith 25 Midway Dr 02134	617 298-8713
17 739-2662	Cartor Broadcasting Co 50 Park Pl St 02134	617 367-9931
17 879-0030	Cartor C 31 East C Cam 02141	617 423-0210
17 541-3948	Cartor C 200 Cambridge St 02136	617 225-2020
17 436-1511	C 210 Harvard Av East Boston 02128	617 782-2118
17 569-6119	C 109 Harvard Cam 02138	617 569-1545
800 622-0213	C & M 41 Northgate 02134	617 491-4822
800 569-8782	C & M 41 Northgate 02134	617 524-9558
	Carter F 51 Hibiscus St 02131	617 327-1105
	Faye & Ricky 20 Columbia Av Mt 02136	617 437-7331
	Francis S 134 Temple W Av 02132	617 323-6781
	Franklin & Anne 705 Mt Auburn Cam 02138	617 354-0798
	Fred 41 Harvard St 02136	617 524-3078
	Fred 16 Harvard St Mt 02136	617 698-1343
	G & B 8 Vernon Ave 02134	617 436-8906
	G T 27 Franklin St 02145	617 623-7121
	Gayle 25 Franklin St 02134	617 825-8322
	Geo S 115 Mass Hill Rd Mt 02138	617 522-3215
	George 125 Boston St 02134	617 367-9548
	Carter Hillside Assoc 107 S Street St 02111	617 456-1689
	Carter Harry P 30 Burns Rd W Av 02132	617 325-5465
	Carter Hide Co Inc 145 Essex St 02131	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 Cam 15 Chatham St 02144	617 323-7990
	J 45 Harvard St 02146	617 730-9483
	J 775 The Pine Wy Roxbury 02132	617 323-5374
	Carter J Jacques MD 1 Crockett Pl Br 02144	617 735-8787
	Carter J M 3410 Columbia Rd S Cam 02138	617 464-1040
	Carter J M Ornamental Ironworks 400 Franklin St 02111	617 436-5353
	Carter J Neal Co 40 Newbury St 02138	617 442-1775
	Carter James 1573 Cambridge St Cam 02138	617 492-1214
	James 412 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
	Carter J A 1200 Cambridge St 02136	617 564-0435
	John 11 Mansfield St 02134	617 426-9094
	John 207 Summer St 02135	617 987-2163
	John 40 Harvard St 02138	617 423-4334
	John 40 Harvard St 02138	617 282-1235
	John D 129 A Summit Av Br 02133	617 734-6199
	K 29 Harvard St 02134	617 265-8656
	K 71 Harvard St 02134	617 282-1593
	K 71 Harvard St 02134	617 282-1593
	Carter Nella E 323 Main St Br 02115	617 267-6483
	Nicholas S F 115 Randolph Av Mt 02136	617 698-5307
	Nick 21 Furlong St 02114	617 267-5222
	Nick & Debbi 136 Harvard Rd Newton 02459	617 527-0480
	Norman G 38 Chickadee Dr 02125	617 822-1203
	P 40 Cranston Pl Br 02135	617 427-4754
	P E 501 E South St Box 02137	617 268-4213
	P L 44 Hutchings Box 02131	617 427-9170
	P R 91 Boyer Cam 02138	617 968-8692
	Paul & Constance 114 Beacon Av W Mt 02133	617 325-3034
	Paul E 501 E South St Box 02137	617 268-4546
	Paul M 27 Union St 02135	617 787-2115
	Carter Pile Driving Inc 27 Beaver Ct Franklin 02102	617 266-5488
	Carter Prudence 40 Franklin Waterbury 02172	617 393-3782
	Prudence 40 Franklin Waterbury 02172	617 926-7063
	Reginald 100 Broadway Cambridge 02142	617 541-2843
	Renee & Andrew 30 Walnut St 02138	617 720-3765
	Carter Rice Donald Building Division Publishing 163 Main Wilmington 01887 Toll Free 800 7 8 7 8 Carl Eric Industrial Prod 613 Main Wilmington Toll Free 800 7 8 7 8 Toll Free 800 7 8 7 8 Toll Free 800 7 8 7 8 Rogers 413 Main Wilmington 01887 978 988-7447 Ingalls Crane 163 Main Wilmington 01887 800 638-1673	
	Carter Richard 2075 Cambridge Av Brighton 02135	617 982-0836
	Richard A 97 Mt Vernon St 02136	617 566-7293
	Carter Richard A MD 1200 Cambridge St 02136	617 267-0710
	Carter Richard K 123 Mount St Box 02137	617 268-0468
	Robert L 175 Rockwood Av Cam 02141	617 864-1535
	Robert L 175 Rockwood Av Cam 02141	617 824-6148
	Royce 18 Salisbury Cha 02129	617 491-6115
	Royce 18 Salisbury Cha 02129	617 241-9418



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 Ln Ipswich 01938	978 356-9960
81 447-4101	Cartagena Lydia 28 Sweet Briar Rd 02331	617 323-7639
90 257-9961	Cartagena Avish F Beach Rd 02319	617 442-9780
17 566-1282	B Had 02334	617 361-5253
17 364-5188	Justicia 50 Decatur Cha 02129	617 241-0152
361-0380	Luzmila 124 Harvard Can 02135	617 491-5621
17 566-4548	M 95 Howe Box 02334	617 323-9713
17 628-8248	Melvin 503 Green Can 02139	617 576-1061
17 445-5116	Carte Nicholas 18 Appleton Boston 02114	617 695-6996
17 822-2962	Cartagena O 4 Bradford Box 02133	617 338-0219
17 427-5712	Carten Thos J Sr & Claire 1 Furlow Ln Mt 02136	617 698-6163
17 569-2698	Carte A 200 Pines Av Cambridge 02142	617 492-4174
17 667-5190	A 200 Pines Av Cambridge 02142	617 492-4174
17 569-1417	Adams 301 Carter St Mt 02136	617 698-7074
17 338-1101	Alice 40 Market Cambridge 02139	617 945-2711
17 825-1193	Andrew F 42 West St 02135	617 625-7623
17 296-1293	Carte Anne MD 1101 Beacon Bn 02444	617 739-1022
17 670-2078	B E 18 Gladstone Av Mt 02136	617 536-6229
17 621-9001	Carte Barbara L MD Turks-New England Medical Center Box 02111	617 296-6911
17 296-4725	Carte Becky Box 02134	617 636-0951
17 542-1521	Bernard J 301 Ashdown E Mt 02136	617 523-4368
17 364-5232	Bibb 25 Midway Dr 02134	617 567-9430
17 541-5649	Bill 301 Ashdown E Mt 02136	617 298-8713
17 739-2662	Carte Broadcasting Co 50 Park Pl Box 02134	617 367-9931
17 879-0030	Carte C 200 Cass 02451	617 423-0210
17 541-3948	Carte C 200 Cass 02451	617 225-0200
17 436-1511	C 210 Townsend Av East Boston 02128	617 782-2118
17 569-4119	C 109 Herman Can 02136	617 569-1545
809 569-8782	C & M 41 Northgate Jan 02134	617 491-8822
	C & M 41 Northgate Jan 02134	617 524-9558
	Carter F 514 Hicks Box 02131	617 327-1105
	Faye & Ricky 20 Columbia Av Box 02136	617 437-7331
	Francis S 134 Temple W Av 02132	617 323-6781
	Franklin & Anne 701 Mt Auburn Can 02138	617 354-0798
	Fred 41 Howard Jan 02136	617 524-3078
	Fred 76 Howley Av Mt 02136	617 698-1343
	G & B 8 Yorker Box 02134	617 436-8906
	G T 27 Fossil Av East 02145	617 623-7121
	Gayle 25 Franklin St 02134	617 823-8322
	Geo S 115 Mount Mt Jan 02136	617 522-3215
	George 52 Madison Box 02134	617 367-9548
	Carter Hillside Assoc 107 S Street Box 02111	617 456-1689
	Carter Harry F 30 Bayview Rd W Av 02132	617 325-5465
	Carter Hide Co Inc 140 Boston St Mt 02136	617 542-7987
	Carter Hilary 41 Harvey Can 02148	617 876-2750
	Horace 301 Walnut Av Rosbury 02138	617 442-5307
	Howard Jr 28 New One Box 02118	617 445-5552
	J Can 15 Chatham Bn 02444	617 232-7990
	J 538 Harvard St 02444	617 730-9483
	J 775 The Pines West Rosbury 02132	617 323-5274
	Carter J Jacques MD 1 Brockton Pl Bn 02444	617 735-8787
	Carter J M 3410 Columbia Rd S Box 02137	617 464-1040
	Carter J M Ornamental Ironworks 100 Franklin Falls 02137	617 436-5353
	Carter J Veal Co 40 Newmarket Box 02138	617 442-1775
	Carle James 1573 Cambridge St Can 02136	617 492-1214
	James 422 Foster Av Rosbury 02138	617 739-2193
	James 31 East Star Rd Cambridge 02141	617 876-8841
	Jane 14 Rosbury Rd Mt 02136	617 361-0773
	Janis 14 Rosbury Rd Mt 02136	617 364-0435
	John 11 Mansfield Bn 02134	617 426-9094
	John 207 Summer St 02135	617 987-2163
	John 40 Woodland St 02135	617 423-4334
	John D 129 A Summit Av Bn 02131	617 282-1235
	J 29 Woodland St 02134	617 265-8656
	K 17 Exposed Dr 02132	617 282-1593
	K 17 Exposed Dr 02132	617 282-1593
	Carter Nellie E 323 Marchessault Box 02115	617 267-6483
	Nicholas S F 115 Randolph Box 02136	617 698-5307
	Nick 21 Furlow Box 02114	617 267-5222
	Nick & Debbi 136 Hermit Rd Newton 02459	617 527-0480
	Norman G 38 Chickadee Dr 02125	617 822-1203
	P 41 Eastwood Pl Box 02135	617 427-4754
	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 02138	617 968-8692
	Paul & Constance 114 Beacon Av W Mt 02133	617 325-3034
	Paul E 501 E South S Box 02137	617 268-4546
	Paul M 27 Union St 02135	617 787-2115
	Carter Pile Driving Inc 27 Beacon St Frankenm 02102	Wellesley Tpk 781.235-0488
	Carter Prudence 40 Franklin Waterbury 02172	617 393-3782
	Prudence 40 Franklin Waterbury 02172	617 926-7063
	Reginald 100 Brookmarch Dr 02124	617 541-2843
	Renée & Andrew 30 Walnut Box 02138	617 720-3765
	Carter Rice David Building Division 163 Main Wilmington 01887 Toll Free 241 7 8 7 Thru.....	800 638-1671
	Carter Richard A Toll Free 241 7 8 7 Thru.....	800 619-7447
	Carter Richard A Toll Free 241 7 8 7 Thru.....	800 648-7447
	Carter Richard A Ingenia Crane 163 Main Wilmington 01887 Toll Free 241 7 8 7 Thru.....	800 638-1673
	Carter Richard 2079 Carleton Av Brighton 02111	617 987-0836
	Carter Richard A MD 1701 Waverley St 02136	617 566-7293
	Carter Richard A 1200 Cambridge St 02134	617 267-0710
	Carter Richard K 123 Merwin St Box 02137	617 268-0468
	Robert L 175 Rockwood Av Can 02141	617 864-1535
	Royce 110 South St 02131	617 424-6148
	Royce 18 Sandway Cha 02129	617 491-6115
	Royce 18 Sandway Cha 02129	617 241-9418



\approx We develop a (conceptual) geography of clusterings

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)

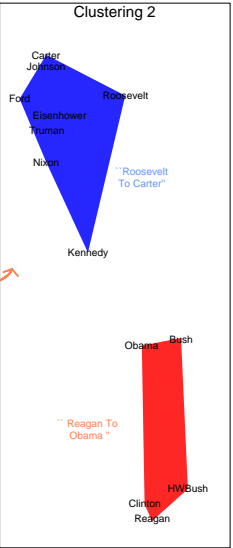
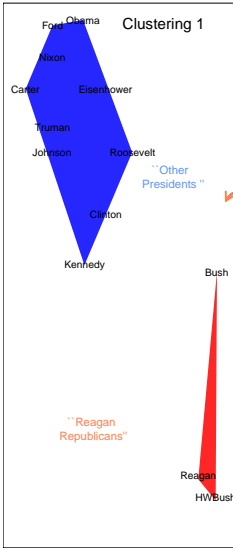
A New Strategy

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 Performance

Evaluating Performance

- Goals:

Evaluating Performance

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

Evaluating Performance

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

Evaluating Performance

- 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 Performance

- 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 Performance

- 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
 - Cluster Quality \Rightarrow RA coders

Evaluating Performance

- 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
 - Cluster Quality \Rightarrow RA coders
 - Informative discoveries \Rightarrow Experienced scholars analyzing texts

Evaluating Performance

- 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
 - Cluster 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)

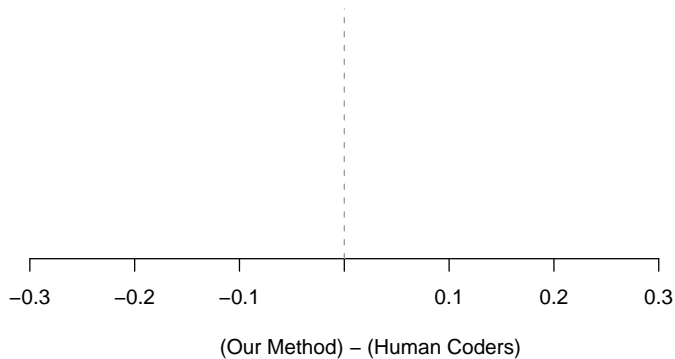
Evaluation 1: Cluster Quality

- **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**

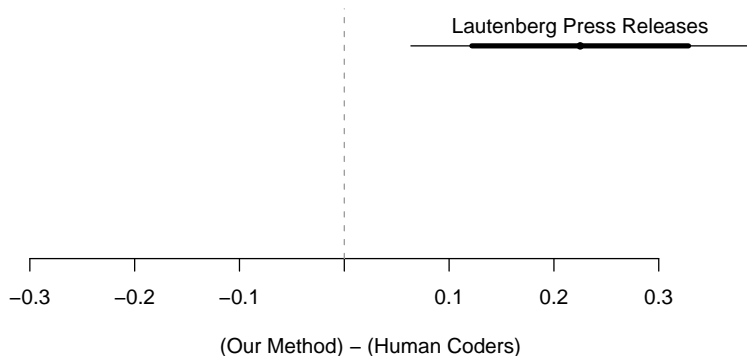
Evaluation 1: Cluster Quality

- **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

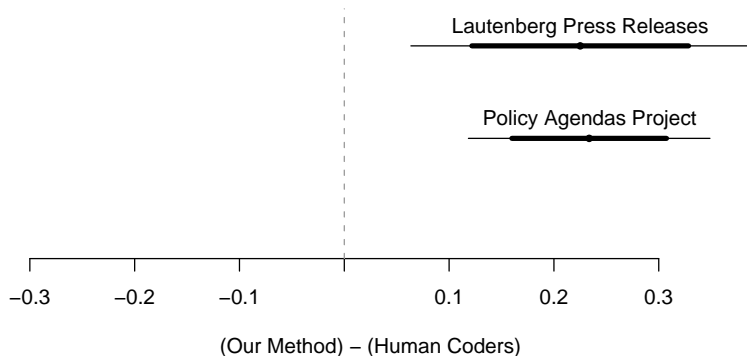


Evaluation 1: Cluster Quality



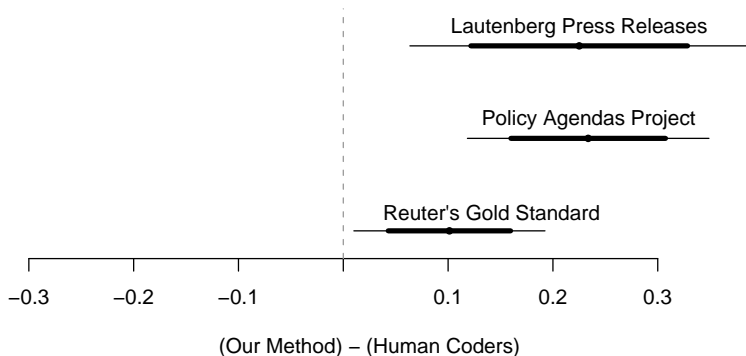
Lautenberg: 200 Senate Press Releases (appropriations, economy, education, tax, veterans, ...)

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

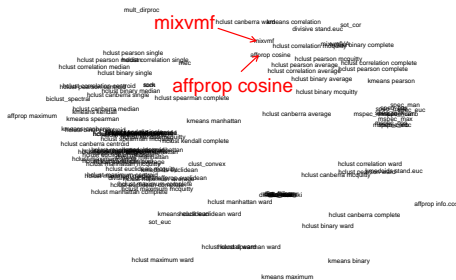
Evaluation 3: What Do Members of Congress Do?

- 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

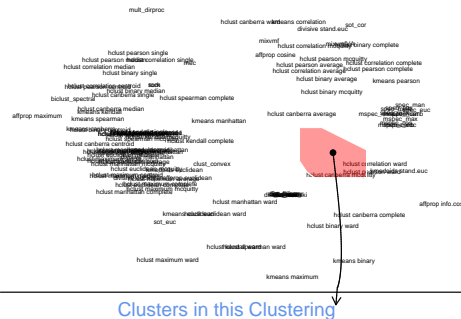


Red point: a **clustering** by
Affinity Propagation-Cosine
(Dueck and Frey 2007)

Close to:

Mixture of von Mises-Fisher
distributions (Banerjee et. al.
2005)

Example Discovery

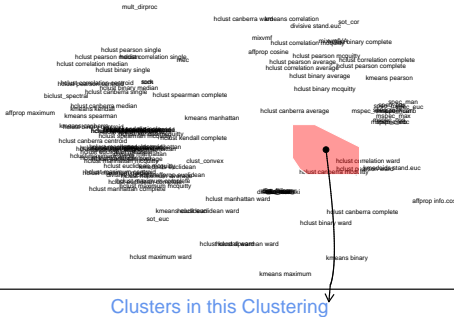


Credit Claiming
Pork

Credit Claiming, Pork:

“Sens. Frank R. Lautenberg (D-NJ) and Robert Menendez (D-NJ) announced that the U.S. Department of Commerce has awarded a \$100,000 grant to the South Jersey Economic Development District”

Example Discovery

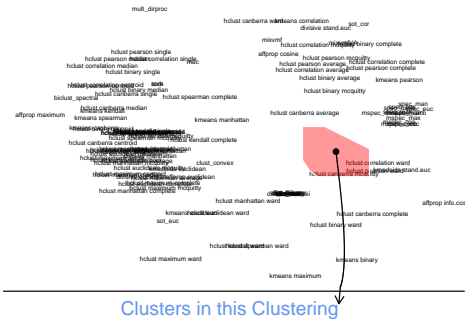


Credit Claiming, Legislation:
“As the Senate begins its recess, Senator Frank Lautenberg today pointed to a string of victories in Congress on his legislative agenda during this work period”

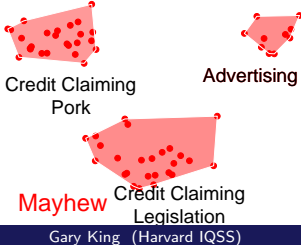
Credit Claiming
Pork

Mayhew
Credit Claiming
Legislation
Gary King (Harvard IQSS)

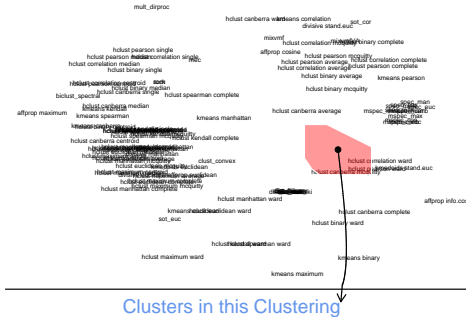
Example Discovery



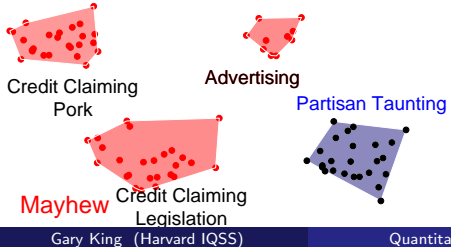
Advertising:
 "Senate Adopts
 Lautenberg/Menendez Resolution
 Honoring Spelling Bee Champion
 from New Jersey"



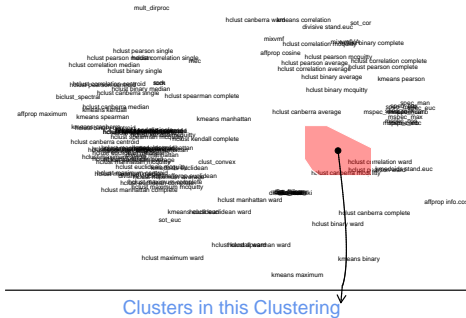
Example Discovery: Partisan Taunting



Partisan Taunting:
“Republicans Selling Out Nation
on Chemical Plant Security”



Example Discovery: Partisan Taunting



Credit Claiming
Pork

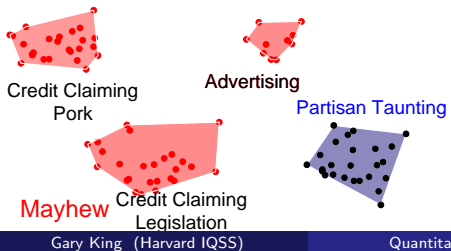
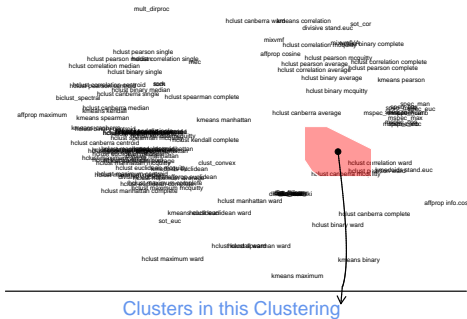
Advertising

Partisan Taunting

Mayhew
Credit Claiming
Legislation
Gary King (Harvard IQSS)

Partisan Taunting:
“Senator Lautenberg’s amendment would change the name of . . . the Republican bill. . . to ‘More Tax Breaks for the Rich and More Debt for Our Grandchildren Deficit Expansion Reconciliation Act of 2006’”

Example Discovery: Partisan Taunting

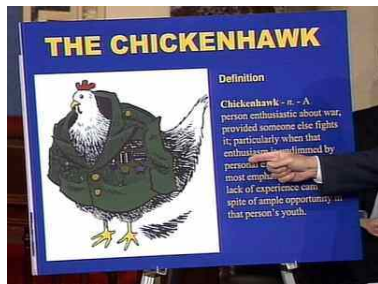


Definition: Explicit, public, and negative attacks on another political party or its members

Taunting ruins deliberation

In Sample Illustration of Partisan Taunting

Taunting ruins deliberation

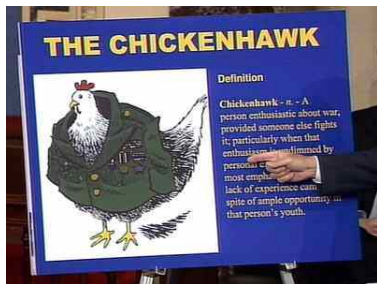


Sen. Lautenberg
on Senate Floor
4/29/04

- "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

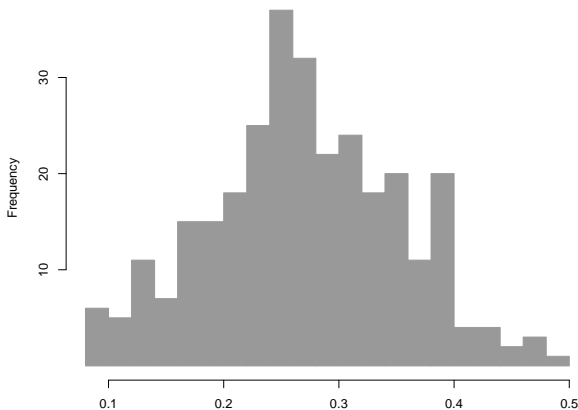
- Discovered using 200 press releases; 1 senator.
- Confirmed using 64,033 press releases; 301 senator-years.

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

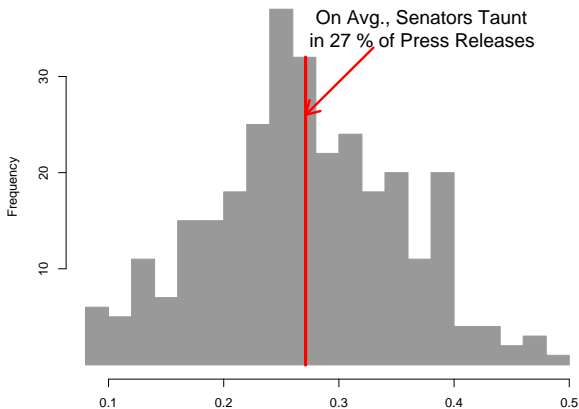
Out of Sample Confirmation of Partisan Taunting

- 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

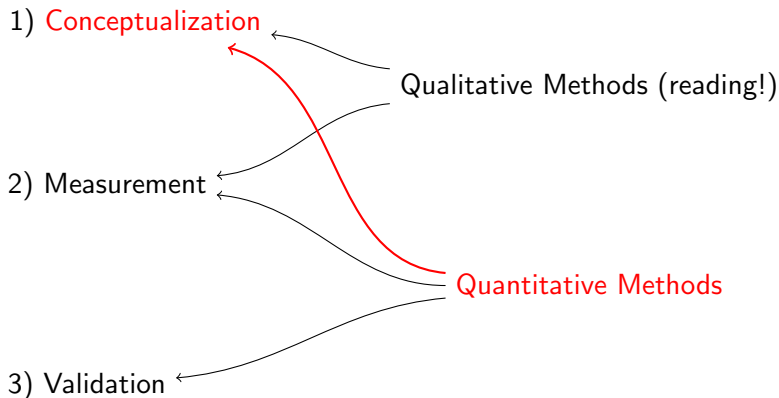


Out of Sample Confirmation of Partisan Taunting

- 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

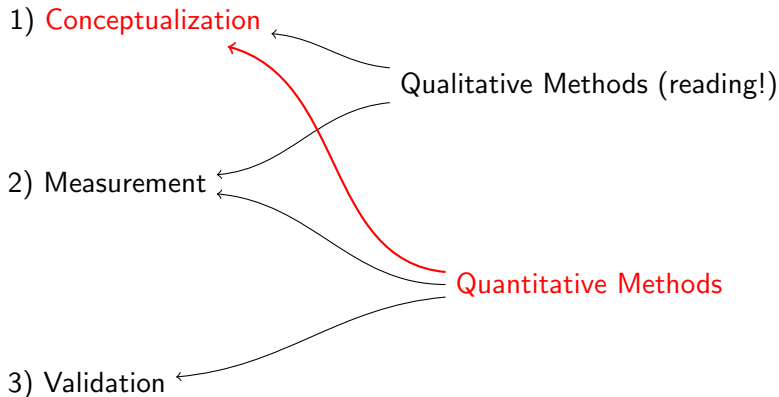


Advancing the Objective of Discovery



Quantitative methods for conceptualization: aiding **discovery**

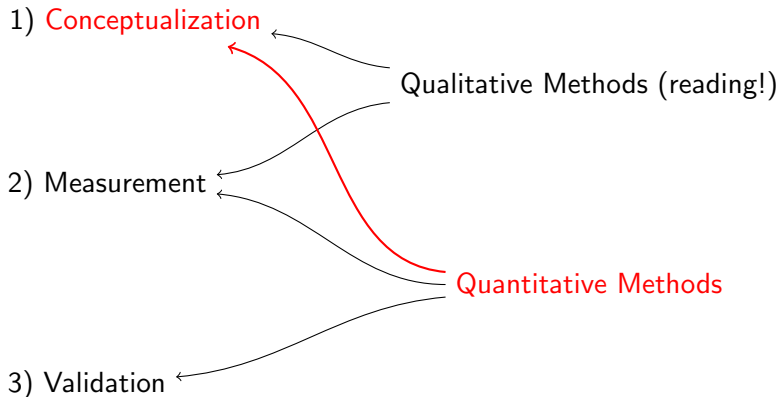
Advancing the Objective of Discovery



Quantitative methods for conceptualization: aiding **discovery**

- Few formal methods designed explicitly for conceptualization

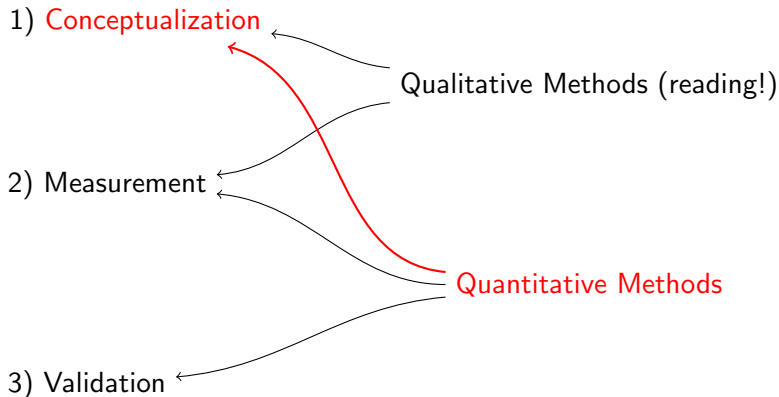
Advancing the Objective of Discovery



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

<http://GKing.Harvard.edu>