

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 Nellie 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 Hed 02134..... 617 361-5253	Franklin & Anne 201 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 W 160 Valley Rd Mill 02186..... 617 698-1343	Norman G 38 Chickawhatch Dr 02125..... 617 822-1201	
M 90 Howe Box 02132..... 617 323-9713	G & B 8 Vardon Dr 02134..... 617 434-8906	P 40 Cranford Pl Box 02135..... 617 437-4754	
Melvin 503 Green Cam 02129..... 617 576-1061	G T 27 Franklin St 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	Gayle 25 Franklin St 02133..... 617 823-0322	P L 44 Hutchings Box 02131..... 617 427-9170	
Cartier 0 4 Bedford Box 02133..... 617 338-0219	George 125 Madison Box 02134..... 617 367-9548	P R 91 Boyer Jan 02138..... 617 968-8692	
Carten Thos Jr S & Claire 17 Franklin St Mill 02135..... 617 698-6163	Carter Hillside Assoc 107 S Street Box 02111..... 617 456-1689	Paul & Constance 114 Franklin St W Box 02131..... 617 325-2036	
17 445-5116	Carter Harry F 30 Bayview Rd W Box 02132..... 617 325-5465	Paul E 501 E South S Box 02137..... 617 268-4546	
17 822-2962	Carter Hide Co Inc 26 Irving St 02133..... 617 542-7987	Paul M 27 Union St 02139..... 617 787-2115	
17 427-5712	A Weber 617 442-5230	Carter Pike Driving Inc 27 Avenue G Frankston 02102..... Wellesley Tpk 781.235-0488	
17 569-2698	Carter Hilary 41 Harvey Cam 02148..... 617 876-2750	Carter Prudence 40 Franklin Waterbury 02127..... 617 393-3782	
17 667-5190	Horac 301 Walnut St Roxbury 02119..... 617 442-5307	Prudence 40 Franklin Waterbury 02127..... 617 926-7063	
17 569-1417	Howard Jr 28 Nona Drive Box 02118..... 617 445-5532	Roginald 106 Brookview Dorchester 02122..... 617 541-8943	
17 338-9110	J Dan 617 354-2658	Renee & Andrew 100 Walnut Box 02118..... 617 720-3765	
17 825-1593	J 31 Chatham Box 02146..... 617 232-7990	Rice Doree Publishing 163 Main Wilmington 01887 Toll Free 800 638-1671	
17 296-1929	J 538 Harvard Box 02146..... 617 730-9483	Carl Sec Industrial Prod 113 Main Wilmington Toll Free 800 619-7447	
17 670-2078	J 775 The Pines West Roxbury 02132..... 617 323-5374	Toll Free 800 719 7 & Thom..... 800 648-7447	
17 621-9001	Jacques M D 1 Brookline Pl Box 02146..... 617 735-8787	Carl..... 978 988-7447	
17 296-4725	Carter J M 3410 Columbia Rd S Box 02137..... 617 464-1040	Carl..... 800 638-1671	
17 542-1521	B E 18 Graduate Ave Mill 02136..... 617 296-6911	Carl..... 800 619-7447	
17 364-5232	Carter Barbara L MD Tufts New England Medical Center Box 02111 Cam..... 617 636-0051	Carl..... 800 648-7447	
17 541-5249	Carter Becky Jo 02134..... 617 523-4368	Carl..... 978 988-7447	
17 739-2662	Bernard J 122 Goodville F Box 02136..... 617 567-3430	Carl..... 800 638-1673	
17 879-0030	B E 18 Graduate Ave Mill 02136..... 617 296-6911	Carter Richard 2079 Cambridge Ave Brighton 02215..... 617 982-0836	
17 436-1511	Bibbiah 25 Midway Dr 02136..... 617 298-8713	Richard A MD 47 Mt Vernon Box 02106..... 617 566-7293	
17 569-4119	Blair 26 Elmwood Box 02108..... 617 367-9031	Carl..... 617 267-0710	
800 569-8782	Carter Broadcasting Co 20 Park Pl Box 02134..... 617 423-0210	Carter Richard K 23 Mather S Box 02127..... 617 268-0408	
	Carter & Bussac Consultants Inc 73 East St Cam 02141..... 617 225-0200	Carl..... 617 864-1535	
	Carter C 200 Commonwealth St 02135..... 617 782-2118	Roger 130 St Bourne Box 02131..... 617 424-6148	
	C 218 Harvard Ave East Boston 02128..... 617 569-1545	Roy 41 Concord Cam 02138..... 617 491-6115	
	C 109 Harvard Cam 02131..... 617 491-8522	Royce 18 Safford Cha 02129..... 617 241-0418	
	C 28 Irving St 02133..... 617 542-4392		
	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 Ln Ipswich 01938	978 356-9960	Carter F. 514 Hickox Ave 02131	617 327-1105	Carter Nella E 323 Marchant Ave Box 02115	617 267-6483
17 447-4101	Cartagena Lydia 28 Sweet Briar 02131	617 323-7639	Faye & Ricky 20 Columbia Ave Box 02136	617 437-7331	Nicholas S F 115 Randolph Ave 02136	617 698-5307
100 257-9961	Cartagena Avish F Beach Rd 02139	617 442-9780	Francis S. 134 Temple W Ave 02132	617 323-6781	Nick & Debbi 215 Fyfield Ave 02116	617 267-5222
17 566-1282	B Had 02136	617 361-5253	Franklin & Anne 705 Mt Auburn Cam 02138	617 354-0798	Norman G 196 Hermit Rd Newton 02459	617 527-0480
17 364-5188	Justica 50 Decatur Cha 02129	617 241-0152	Fred 41 Haverhill Aven 02136	617 524-3078	38 Chickadee Rd Der 02125	617 822-1203
361-0380	Luzmila 124 Harvard Cam 02138	617 491-5621	Fred W. Haverhill Ave 02136	617 698-1343	40 Cranston Rd Box 02115	617 427-4754
17 566-4548	M 90 Howe St 02139	617 323-9713	G & B. 8 Vardon Ave 02134	617 436-8906	P E 501 E South St Box 02137	617 268-4213
17 628-8248	Melvin 503 Green Cam 02139	617 576-1061	G T 27 Fyfield Ave 02136	617 623-7121	P L 44 Hutchings Box 02115	617 427-9170
17 445-5116	Carte Nicholas 18 Appleton Boston 02114	617 695-6996	Gayle 25 Franklin Der 02124	617 823-0322	P R 91 Brewer Ave 02138	617 968-8692
17 822-2962	Cartagena O 4 Bradford Box 02118	617 338-9219	Geo S 115 Mass Hill Rd Box 02138	617 522-3215	Paul & Constance 114 Adams Ave W Box 02110	617 325-3034
17 427-5712	Carten Thos J Sr & Claire 1 Fyfield St Mt 02116	617 698-6163	George 120 Howe Ave 02134	617 367-9548	Paul M 207 Union St 02139	617 787-2115
17 569-2698	Carte Thos & Kathleen 50 Thompson Ln Mt 02136	617 696-6919	Carter Harry F 107 S Street Box 02111	617 456-1689	Paul M 207 Union St 02139	617 787-2115
17 667-5190	Carte A A 200 Pioneer Av Cambridge 02142	617 492-4174	Carter Hide Co Inc 140 Bunker Hill Rd W Box 02112	617 325-5465	Prangman 02102	Wellesley Tpk-781.235-0488
17 569-1417	Adams 301 Carter St Mt 02136	617 698-9074	Carter Hilary 41 Harvey Cam 02148	617 876-2750	Carter Prudence 40 Franklin Waterbury 02172	617 393-3782
17 338-9110	Alice 40 Market Cambridge 02139	617 945-2711	Horace 301 Walnut Av Roxbury 02119	617 442-5307	Prudence 40 Franklin Waterbury 02172	617 926-7063
17 825-9195	Andrew F 42 Mt St 02135	617 625-7623	Howard J 301 Walnut Av Roxbury 02119	617 442-5307	Reginald 100 Brookside Cambridge 02142	617 541-2843
17 296-1293	Carter Anne MD 1101 Beacon Ave 02144	617 739-1022	Howard Jr 28 New One Box 02118	617 445-5532	Renee & Andrew 100 Brookside Cambridge 02142	617 541-2843
17 670-2078	B E 10 Gladstone Ave Mt 02136	617 296-6911	J C 15 Chatham St 02144	617 232-7990	Carter Rice Doan 301 Walnut St 02138	617 720-3765
17 621-9001	Carte Barbara L MD Tufts New England Medical Center Box 02111	617 436-0951	J S 4775 The Pines West Roxbury 02132	617 323-5274	Carl Tud Free-Dad '7 & Th... 800 638-1671	
17 296-4725	Carter Becky 90 02134	617 523-4368	Carter J Jacques MD 1 Crockett Pl Br 02144	617 735-8787	Carl Tud Free-Dad '7 & Th... 800 648-7447	
17 542-1521	Bernard J 371 Newbury Boston 02116	617 536-6229	Carter J M 3410 Columbia Rd S Cam 02138	617 464-1040	Carl Tud Free-Dad '7 & Th... 800 648-7447	
17 364-5232	Bibbath 25 Midway Der 02124	617 298-8713	Carter J M Ornamental Ironworks Pondside Falls 01743	617 436-5353	Carl Tud Free-Dad '7 & Th... 800 648-7447	
17 541-5649	Carte Broadcasting Co 50 Park Pl Box 02136	617 423-0210	Carter J Neal Co 40 Newmarket Rd 02138	617 442-1775	Carl Tud Free-Dad '7 & Th... 800 648-7447	
17 739-2662	Carter Business Consultants Inc 31 East Cam 02141	617 225-0200	Carter J Paul 157 Cambridge St Cam 02138	617 492-1214	Carl Tud Free-Dad '7 & Th... 800 648-7447	
17 879-0030	Carter C 2000 Gesserts Av Br 02135	617 782-2118	James 452 Foster Ave Roxbury 02119	617 739-2193	Carl Tud Free-Dad '7 & Th... 800 648-7447	
17 436-1511	C 210 Harvard Av East Boston 02128	617 569-1545	James 31 East Star Rd Cambridge 02141	617 876-8841	Carl Tud Free-Dad '7 & Th... 800 648-7447	
17 569-4119	C 109 Harvard Cam 02138	617 491-4822	Jane L 34 Rosbury Rd Mt 02146	617 361-0773	Carl Tud Free-Dad '7 & Th... 800 648-7447	
800 669-8782	C & M 41 Northgate Ave 02134	617 524-9558	Janice 14 Adams Rd Newton 02459	617 564-0435	Carl Tud Free-Dad '7 & Th... 800 648-7447	
			John 11 Mansfield Dr 02134	617 987-2163	Carl Tud Free-Dad '7 & Th... 800 648-7447	
			John 207 Summer St 02125	617 423-4134	Carl Tud Free-Dad '7 & Th... 800 648-7447	
			John 40 Westfield Der 02125	617 282-1235	Carl Tud Free-Dad '7 & Th... 800 648-7447	
			June O 129 A Summit Av Br 02133	617 734-6109	Carl Tud Free-Dad '7 & Th... 800 648-7447	
			K 179 Inwood Cambridge 02142	617 265-8656	Carl Tud Free-Dad '7 & Th... 800 648-7447	
			K 17 Exposed Der 02117	617 282-1593	Carl Tud Free-Dad '7 & Th... 800 648-7447	



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	
17 447-4101	Cartagena Lydia 28 Sweet Briar Rd 02131	617 323-7639	
90 257-9961	Cartagena Avish F Beach Rd 02139	617 442-9780	
17 566-1282	B Had 02136	617 361-5253	
17 364-5188	Lucille 174 Harvard Can 02136	617 491-5621	
361-0380	M 95 Howe Box 02136	617 323-9713	
17 566-4548	Melvin 503 Green Can 02139	617 576-1061	
17 628-8248	Carte Nicholas 18 Appleton Boston 02114	617 695-6996	
17 445-5116	Carlton 4 414ford Box 02138	617 338-0219	
17 822-2962	Carton 50 Thompson Ln Mt 02136	617 696-6919	
17 427-5712	A Heber A 200 Pitman Av Cambridge 02142	617 492-4174	
17 569-2698	A 21 Beetham Wy Haverhill 02119	617 442-1219	
17 667-5190	A M 250 Main St Av Box 02115	617 266-7153	
17 569-1417	Adams 361 Carter St Mt 02136	617 698-9074	
17 338-1107	Alice 40 Market Cambridge 02139	617 945-2711	
17 822-1959	Artner 62 Mt Airy St Box 02138	617 625-7623	
17 296-1193	Arter Anne MD 1161 Beacon Bn 02144	617 739-1022	
17 670-2078	B E 18 Gladstone Av Mt 02136	617 536-6229	
17 621-9001	Barber Barbara L MD Tufts New England Medical Center Box 02111	617 296-6911	
17 296-4725	Barber Becky MD 02114	617 636-0951	
17 542-1521	Bernard J 371 Newbury Boston 02116	617 523-4368	
17 364-5232	Bibb 25 Midway Dr 02136	617 567-9430	
17 541-5649	Biggs 18 W Newbury St 02138	617 298-8713	
17 739-2662	Carter Broadcasting Co 50 Park Pl Box 02116	617 367-9931	
17 879-0030	Carter C 2000 Cavendish Av St 02135	617 423-0210	
17 541-3948	C 210 Fremont Av East Boston 02128	617 225-0200	
17 436-1511	C 109 Harvard Can 02136	617 782-2118	
17 569-4119	C 8 111 Cambridge St 02139	617 569-1545	
909 569-8782	C & M 41 Northgate Jct 02134	617 491-8822	
		617 524-9558	
		617 282-1193	
		617 327-1105	
		617 437-7331	
		617 323-6781	
		617 354-0798	
		617 524-3078	
		617 698-1343	
		617 436-8906	
		617 623-7121	
		617 823-8322	
		617 522-3215	
		617 367-9548	
		617 456-1689	
		617 325-5465	
		617 542-7987	
		617 876-2750	
		617 442-5307	
		617 445-5552	
		617 354-2658	
		617 232-7990	
		617 730-9483	
		617 323-5274	
		617 735-8787	
		617 464-1040	
		617 436-5353	
		617 442-1775	
		617 492-1214	
		617 739-2193	
		617 876-8841	
		617 364-0773	
		617 964-0435	
		617 426-9094	
		617 987-2163	
		617 423-4334	
		617 282-1235	
		617 734-6109	
		617 265-8656	
		617 282-1593	
		617 267-6483	
		617 698-5307	
		617 698-0713	
		617 327-0480	
		617 627-4754	
		617 268-8213	
		617 427-9170	
		617 968-8692	
		617 325-3034	
		617 268-4546	
		617 787-2115	
		617 268-4546	
		617 393-3782	
		617 926-7063	
		617 541-2843	
		617 720-3765	
		617 838-1671	
		617 987-7447	
		617 608-7447	
		617 988-7447	
		617 638-1673	
		617 987-0836	
		617 566-7293	
		617 267-0710	
		617 268-0448	
		617 864-1535	
		617 424-6148	
		617 491-6115	
		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)

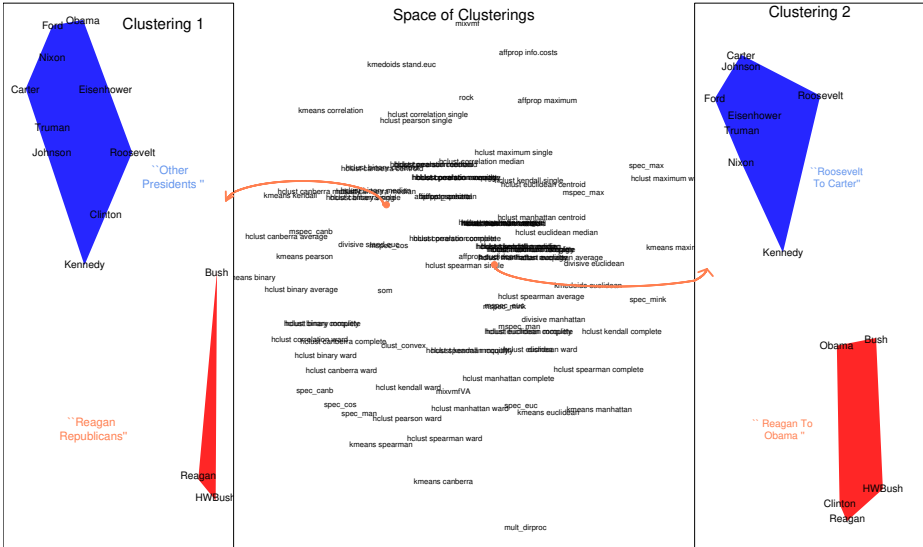
A New Strategy

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)

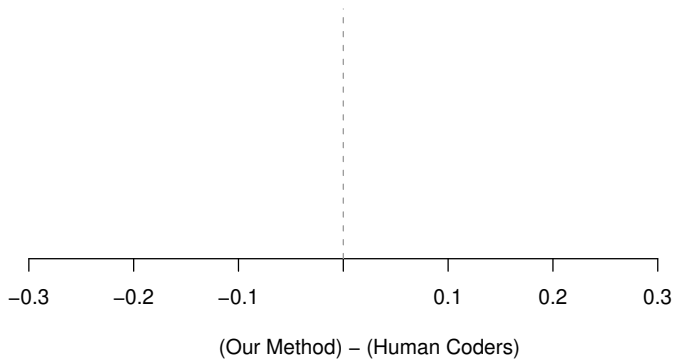
Evaluation 1: Cluster Quality

- **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**

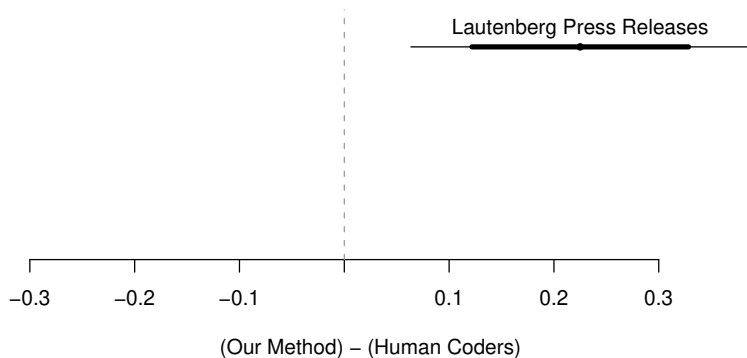
Evaluation 1: Cluster Quality

- **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

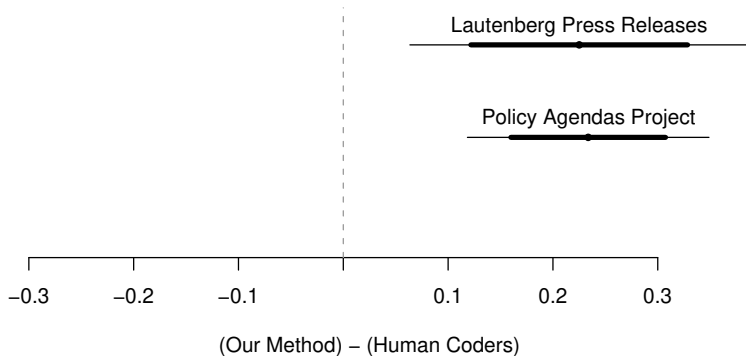


Evaluation 1: Cluster Quality



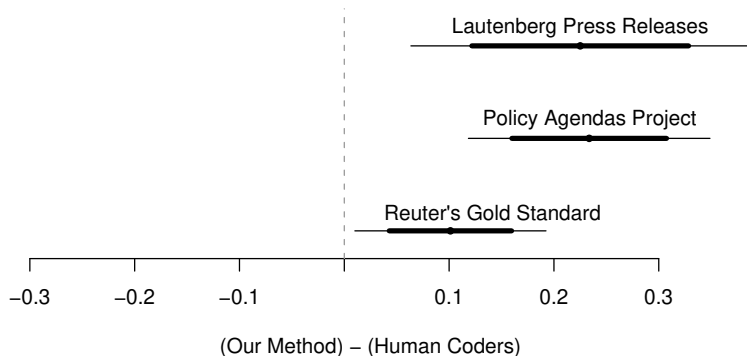
Lautenberg: 200 Senate Press Releases (appropriations, economy, education, tax, veterans, ...)

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

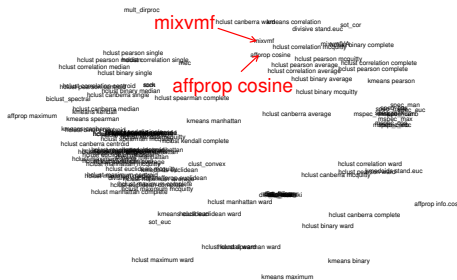
Evaluation 3: What Do Members of Congress Do?

- 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

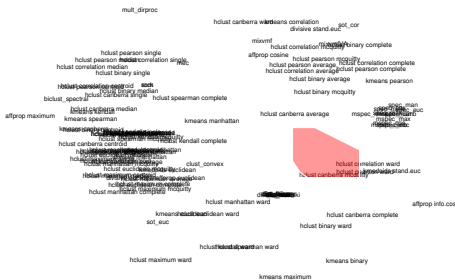


Red point: a **clustering** by Affinity Propagation-Cosine (Dueck and Frey 2007)

Close to:

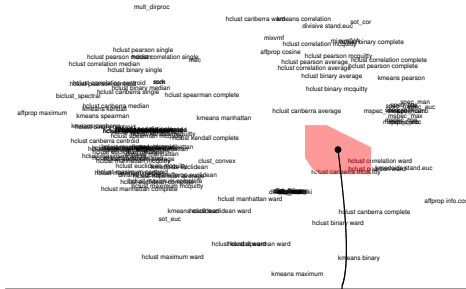
Mixture of von Mises-Fisher distributions (Banerjee et. al. 2005)

Example Discovery



Found a **region** with particularly insightful clusterings

Example Discovery



Clusters in this Clustering



Credit Claiming
Pork



Advertising



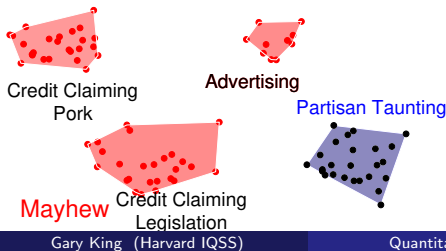
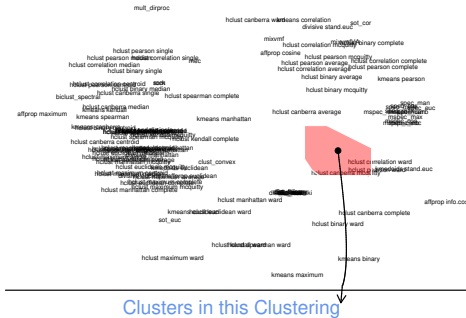
Mayhew
Credit Claiming
Legislation

Gary King (Harvard IQSS)

Advertising:

“Senate Adopts
Lautenberg/Menendez Resolution
Honoring Spelling Bee Champion
from New Jersey”

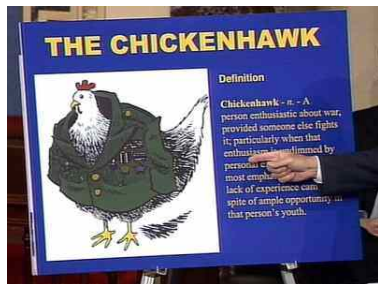
Example Discovery: Partisan Taunting



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’”

In Sample Illustration of Partisan Taunting

Taunting ruins deliberation

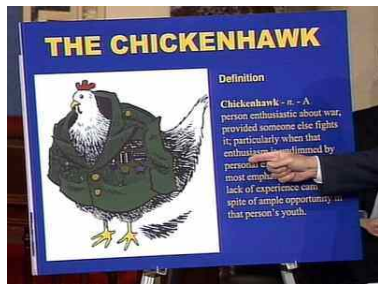


Sen. Lautenberg
on Senate Floor
4/29/04

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

In Sample Illustration of Partisan Taunting

Taunting ruins deliberation

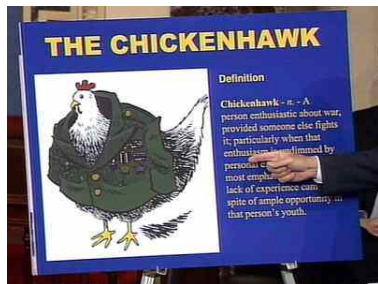


Sen. Lautenberg
on Senate Floor
4/29/04

- "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]

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]
- "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

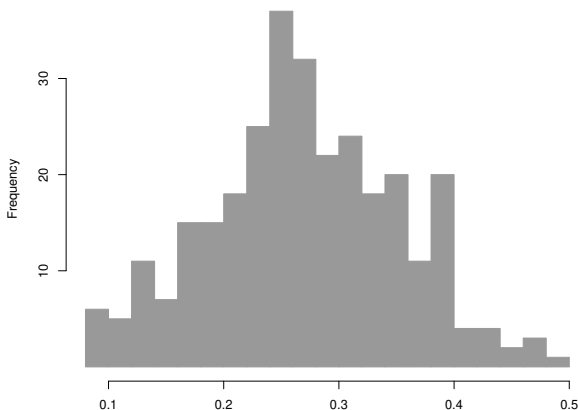
- Discovered using 200 press releases; 1 senator.
- Confirmed using 64,033 press releases; 301 senator-years.

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

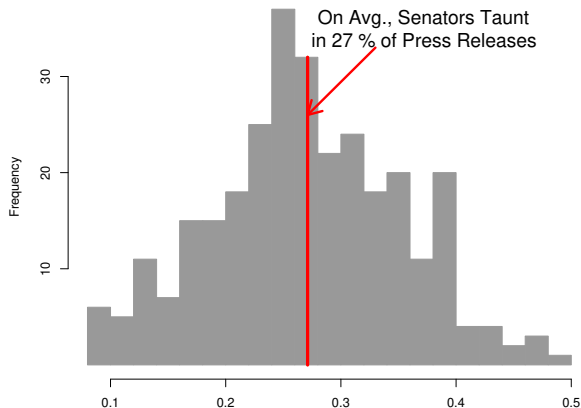
Out of Sample Confirmation of Partisan Taunting

- 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

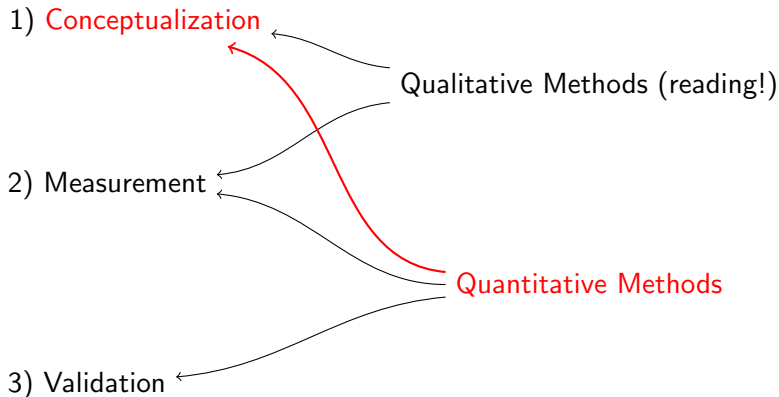


Out of Sample Confirmation of Partisan Taunting

- 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

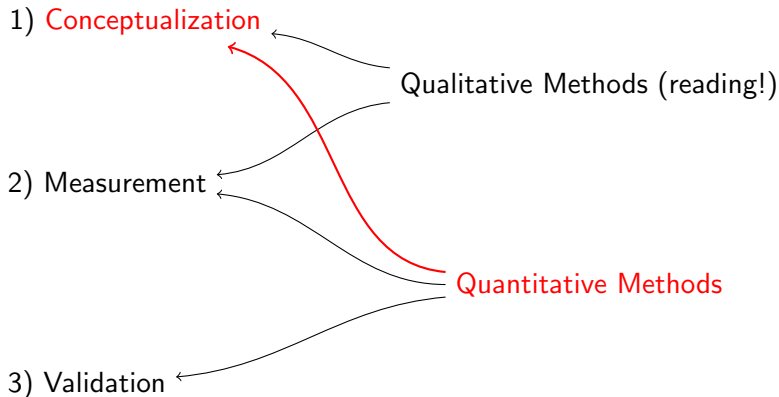


Advancing the Objective of Discovery



Quantitative methods for conceptualization: aiding **discovery**

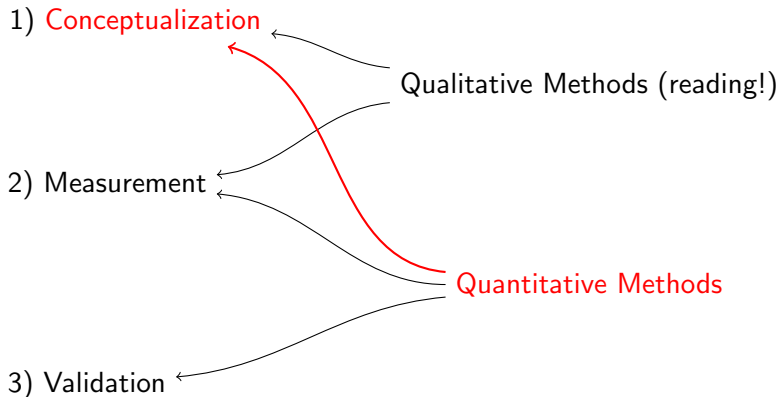
Advancing the Objective of Discovery



Quantitative methods for conceptualization: aiding **discovery**

- Few formal methods designed explicitly for conceptualization

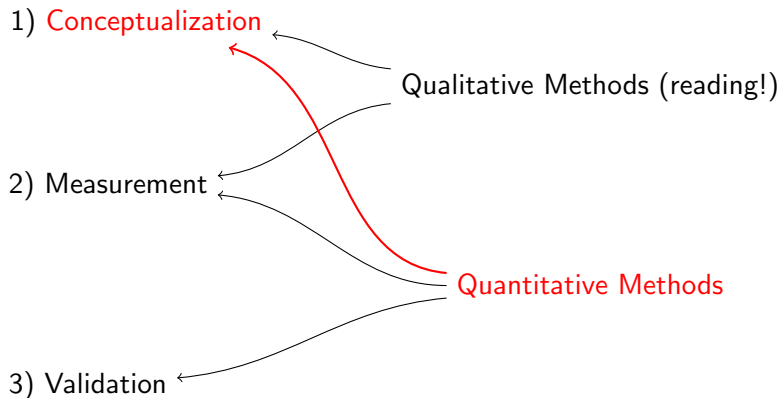
Advancing the Objective of Discovery



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>