

Advanced Quantitative Research Methodology, Lecture Notes: Text Analysis II: Unsupervised Learning via Cluster Analysis¹

Gary King

Institute for Quantitative Social Science
Harvard University

¹© Copyright 2012 Gary King, All Rights Reserved.

A Method for Computer Assisted Conceptualization

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).
- **Cluster Analysis**: simultaneously (1) invents categories and (2) assigns 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).
- **Cluster Analysis**: simultaneously (1) invents categories and (2) assigns documents to categories
- We focus on unstructured text; methods apply more broadly.

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).
- **Cluster Analysis**: simultaneously (1) invents categories and (2) assigns documents to categories
- We focus on unstructured text; methods apply more broadly.
- Main goal: Switch from **Fully Automated** to **Computer Assisted**

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**
 - **Easy in theory:** list all 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**
 - **Easy in theory:** list all clusterings; choose the best
 - **Impossible in practice:** Too hard for us 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**
 - **Easy in theory:** list all clusterings; choose the best
 - **Impossible in practice:** Too hard for us 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**
 - **Easy in theory:** list all clusterings; choose the best
 - **Impossible in practice:** Too hard for us 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**
 - **Easy in theory:** list all clusterings; choose the best
 - **Impossible in practice:** Too hard for us 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**
 - **Easy in theory:** list all clusterings; choose the best
 - **Impossible in practice:** Too hard for us 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
- **Question: How to organize clusterings so humans can understand?**

Our Idea: Meaning Through Geography

Set of clusterings

Our Idea: Mapping Through Geography

Set of clusterings \approx
A list of unconnected addresses

wide at SuperPages.com

		195	Car	C
37 566-1282	Cartage New England Inc	978 356-9960		
81 447-4101	Cartagena Lydia			
90 257-9961	Cartagena Aveth			
361-0380	Carte Nicholas			
37 566-4548	Carter Nicholas			
37 628-8248	Carter Thos J Sr & Claire			
37 445-5116	Thomas & Kathleen			
37 822-9902	Carter A Inc			
37 427-5712	A Harbor			
37 569-2698	A 31 Seabrook Wy			
37 667-5190	A 200 Putnam Ave			
37 569-1417	Adams 381 Carter St			
37 338-0110	Alice 100 Kilmackree Rd			
37 825-9195	Allice 40 Market Cambridge			
37 296-1593	Allice 42 West St			
37 670-2078	Allice 1100 Essex St			
37 623-9001	Allice 1100 Essex St			
37 296-4725	Allice 1100 Essex St			
37 542-1521	Bernard J			
37 364-5232	Bithiah 25 Melrose Dr			
37 541-5649	Bithiah 25 Melrose Dr			
37 739-2662	Carter Broadcasting Co			
37 879-0030	Carter C 2000 Cambridge St			
37 936-1511	C 2000 Cambridge St			
37 569-4119	C 2000 Cambridge St			
37 569-4782	C & M 41 Burroughs Jan			
37 327-1105	Faye & Ricky			
37 437-7331	Francis S 134 Temple W			
37 323-6781	Franklin & Anne			
37 354-0798	Fred 40 Haverhill Jan			
37 524-3078	Fred 40 Haverhill Jan			
37 698-1343	G & R 8 Haverhill Dr			
37 436-8906	G T 27 Franklin St			
37 623-7121	Gayle 25 Franklin St			
37 825-0322	Geo S 115 Main Hill Rd			
37 522-3215	George 105 Madison St			
37 367-9548	Carter Halliday Associate			
37 456-1689	Carter Harry F			
37 325-5465	Carter Hide Co Inc			
37 542-7987	Carter Hilary 41 Harvey Can			
37 876-2750	Horace			
37 442-5307	Howard Jr 28 Neve One Box			
37 445-5552	J Cam			
37 354-2658	J 31 Chatham Ave			
37 232-7990	J 518 Harvard Ave			
37 730-9483	J 775 Wyomissing Wy			
37 323-5374	J 1000 Main St			
37 735-8787	Carter J M			
37 464-1040	3410 Columbia St			
37 436-5353	Carter J M Ornamental Ironworks			
37 442-1775	Carter J Neal Co			
37 492-1214	James 157 Cambridge St			
37 739-2193	James 157 Cambridge St			
37 876-8841	James 31 Gold Star Rd			
37 361-0773	Jan L 34 Broadview Rd			
37 964-0435	Jan L 34 Broadview Rd			
37 436-5994	Jeffrey 40 Warren Ave			
37 987-2163	John 11 Mainfield Rd			
37 423-4334	John 107 Summer Box			
37 282-1535	John 40 Westwood Rd			
37 734-6199	June O 109 A Summit Ave			
37 265-9456	K 200 Wyomissing Wy			
37 282-1593	K 17 Exford Cambridge			
37 267-6483	Carter Nellie E			
37 698-5307	Nicholas S F			
37 267-5222	Nick 21 Fairfield Box			
37 698-0713	Nick & Debbi			
37 527-0480	Nicole 146 Vermont Rd			
37 822-1203	Norman G			
37 437-4754	P 40 Cranston Pl			
37 268-8213	P E 501 E South St			
37 427-9170	P L 40 Haverhill Box			
37 983-8692	Paul & Constance			
37 325-2036	Paul 114 Aspen Ave			
37 268-4546	Paul E 501 E South St			
37 787-2115	Paul M 27 Union St			
37 876-2750	Carter Pile Driving Inc			
37 393-3782	Carter Prudence			
37 926-7063	Prudence 40 Franklin Wadsworth			
37 541-2843	Reginald 100 Broadview Cambridge			
37 720-3765	Renee & Andrew			
800 638-1671	Carter Rice Dowd			
800 619-7447	Carl Rice Dowd			
800 648-7447	Carl Rice Dowd			
978 988-7447	Carl Rice Dowd			
800 638-1673	Carl Rice Dowd			
37 987-0836	Carter Richard			
37 566-7293	Richard A 8748 Vernon St			
37 267-0710	Carter Richard A			
37 268-9448	Richard K 127 Main St			
37 864-1535	Roger 150 St Pauline Box			
37 491-6115	Roy 40 Concord Ave			
37 241-0418	Royce 18 Sumner Dr			

Our Idea: Meaning Through Geography

Set of clusters \approx

A list of unconnected addresses

wide at SuperPages.com

	195	Car	C
Carterge New England Inc			
37 566-1282	26 Allen Ln Ipswich MA 01938	978 356-9960	
Carterge Lydia			
81 447-4101	38 Sweet Rd 02131	617 323-7639	
Carterge Aveth			
90 257-9961	1 Fairview Rd 02139	617 442-9780	
	R 104 02136	617 361-5253	
37 566-1282	Justice 50 Decatur Cha 02129	617 241-0152	
37 564-5188	Lucilla 124 Harvard Cam 02139	617 491-5621	
	M 95 Howe Rd 02133	617 323-9713	
361-0380	Melvin 501 Green Cam 02139	617 576-1061	
Carte Nicholas			
37 566-4548	38 Appleton Boston 02134	617 695-6996	
	Carterge O 4480rd Wm 02131	617 338-9219	
37 628-8248	Carten Thos J Sr & Claire		
	1 Ivesdale Rd MA 02136	617 698-6163	
37 445-5116	Thomas & Kathleen		
	50 Thompson Ln MA 02136	617 696-6919	
37 822-9902	Carter A Sr 02133	617 339-2257	
37 422-5712	A Hubert	617 442-5230	
37 569-2698	A 31 Bethune Wy Rndary 02139	617 442-1219	
	A 200 Putnam Av Cambridge 02139	617 492-4174	
37 667-5190	A M 255 Mchua Av Bos 02135	617 266-7153	
37 569-1417	Adams 361 Carter St MA 02138	617 698-9074	
37 338-1110	Allice 108 Elmwood Rte 02136	617 425-0193	
	Allice 40 Market Cambridge 02139	617 945-2711	
	Andrew F 42 Wal St Som 02143	617 625-7635	
37 825-1919	Carter Anne MD	617 739-1022	
37 296-1593	1161 Beacon Wm 02144		
37 670-2078	B E 108 Graduate Av Mt 02136	617 536-6329	
37 623-9001	Carter Barbara L MD	617 296-6911	
	Tuffs New England Medical Res 02131		
37 296-4725	Carter Becky Sr 02134	617 636-9051	
	Carter Adhene	617 523-4368	
37 542-1521	Bernard J		
	22 Cambridge St Bos 02136	617 567-3430	
37 364-5232	Bithiah 25 Midway Rte 02124	617 298-8713	
37 541-5649	Bliss 312 Newbury St 02138	617 367-9931	
37 739-2662	Carter Broadcasting Co		
	26 Park Pl Bos 02134	617 423-0210	
37 879-0030	C 21 2nd St Cam 02141	617 225-0200	
37 541-3948	Carter C 200 Cambridge Av Bos 02135	617 782-2118	
37 936-1511	C 219 Harvard Av East Boston 02128	617 569-1545	
37 569-4119	C 109 Harvard Cam 02138	617 491-4822	
	C 109 Harvard Cam 02138	617 296-4392	
37 569-4782	C & M 41 Burroughs Jan 02136	617 524-5595	
Carter F 24 Hillock Bos 02131			
	617 327-1105		
Faye & Ricky			
	107 Columbia Av Bos 02136	617 437-7331	
	Francis S 134 Temple W Av 02132	617 323-6781	
Franklin & Anne			
	291 Mt Auburn Cam 02138	617 354-0798	
	Fred 42 Harvard Cam 02136	617 524-3078	
	Fred 56 Harvard Av Mt 02138	617 698-1343	
	G & E 8 Wynden Der 02134	617 436-8906	
	G T 27 Fyfield Av Som 02145	617 623-7121	
	Gayle 25 Franklin Der 02134	617 825-0322	
	Geo S 115 Main Mt Rd Jan 02138	617 522-3215	
	George 125 Nashua Bos 02131	617 367-9548	
Carter Halliday Associate			
	107 S Street Bos 02111	617 456-1689	
Carter Harry F			
	26 Baring Jct Rd Wm 02132	617 325-5465	
Carter Hide Co Inc			
	141 Somers Bos 02133	617 542-7987	
	Carter Hilary 41 Harvey Cam 02140	617 876-2750	
Horace			
	301 Watline Av Rndary 02139	617 442-5307	
	Howard Jr 38 Neha One Bos 02118	617 445-5552	
	J Cam	617 354-2658	
	J 35 Chatham Bos 02146	617 232-7990	
	J 35 Harvard Bos 02144	617 730-9483	
	J 775 Wy Windy Woodbury 02135	617 323-5574	
	Carter J Jacques MD		
	1 Broadview Pl Bos 02144	617 735-8787	
Carter J M			
	3410 Columbia Rd S Bos 02137	617 464-1040	
	Carter J M Ornamental Ironworks		
	Prudential Tolls 617 436-5353		
Carter J Veal Co			
	40 Newbury St Bos 02138	617 442-1775	
Carter James			
	1573 Cambridge St Cam 02136	617 492-1214	
	James 302 Fisher Av Rndary 02136	617 739-2193	
	James 112 Oak St Rd Cambridge 02140	617 876-8841	
	Jas L 34 Broadview Rd Mt 02136	617 361-0773	
	James 124 Adams Rd Newton 02458	617 564-0435	
	Jeffrey 41 Warren Av Som 02136	617 426-5994	
	John 11 Mansfield Rd 02134	617 987-2163	
	John 307 Summer Bos 02137	617 423-4334	
	John 40 Watline Rd Der 02125	617 282-1235	
	June O 129 A Summit Av Bos 02138	617 734-6109	
	K 109 Wynden Av Rndary 02136	617 265-9456	
	K 17 Elwood Der 02123	617 282-1593	
Carter Nellie E			
	323 Marchette Av Bos 02135	617 267-6483	
Nicholas S F			
	115 Randolph Av Mt 02136	617 698-5307	
	Nick 21 Fairfield Bos 02114	617 267-5222	
Nick & Debbi			
	156 Vermont Rd Newton 02459	617 527-0480	
Norman G			
	38 Chickadee Der 02125	617 822-1203	
	38 Chickadee Der 02125	617 427-4754	
	P E 501 E South St Bos 02137	617 268-4213	
	P L 44 Matthews Bos 02131	617 427-9170	
	P R 91 Boyer Jan 02134	617 968-8692	
Paul & Constance			
	114 Adams Av Wm 02130	617 325-2036	
	Paul E 501 E South St Bos 02137	617 268-4546	
	Paul M 27 Green Av 02135	617 787-2115	
Carter Pile Driving Inc 17 Beaver Ct			
	Franklin 02102	Wellesley Tolls 781 235-0488	
Carter Prudence			
	40 Franklin Watline 02127	617 393-3782	
Prudence			
	40 Franklin Watline 02127	617 926-7063	
Reginald			
	106 Broadview Cambridge 02215	617 541-2843	
Renee & Andrew			
	10 Walnut Bos 02138	617 720-3765	
Carter Rice Dowd			
	Baker Dennis Publishing 163 Main Wilmington 01887		
	Toll Free 0-811 & Thon	800 638-1671	
	Cost Inc Industrial Prod 413 Main Wilmington		
	Toll Free 0-811 & Thon	800 619-7447	
	Toll Free 0-811 & Thon	800 648-7447	
	Headquarters 413 Main Wilmington 01887		
	Call	978 988-7447	
	Inglis Crane 363 Main Wilmington 01887		
	Call 0-811 & Thon	800 638-1673	
Carter Richard			
	2075 Carver Av Brighton 02137	617 987-0836	
	Richard A 974 Vernon Bos 02136	617 566-7293	
Carter Richard A MD			
	170 Wynden Av Som 02136	617 267-0710	
Carter Richard K			
	133 Merwin St 02137	617 268-9448	
	Robert L 175 Rockdale Av Cam 02141	617 864-1535	
	Roger 150 St Pauls Bos 02131	617 424-6148	
	Roy 41 Concord Rd 02138	617 491-6115	
	Royce 185 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
37 566-1282	Cartage New England Inc	157 356-9960	
81 447-4101	Cartagena Lydia	157 356-9960	
90 257-9961	Cartagena Avelis	157 356-9960	
37 566-1282	Cartagena Avelis	157 356-9960	
37 566-1282	Cartagena Avelis	157 356-9960	
361-0380	Carte Nicholas	157 356-9960	
37 566-4548	Carte Nicholas	157 356-9960	
37 628-8248	Carten Thos J Sr & Claire	157 356-9960	
37 445-5116	Carten Thos J Sr & Claire	157 356-9960	
37 822-9992	Carte A Inc	157 356-9960	
37 422-5712	Carte A Inc	157 356-9960	
37 569-2698	Carte A Inc	157 356-9960	
37 667-5190	Carte A Inc	157 356-9960	
37 569-1417	Carte A Inc	157 356-9960	
37 822-9992	Carte A Inc	157 356-9960	
37 296-1593	Carte A Inc	157 356-9960	
37 670-2078	Carte A Inc	157 356-9960	
37 623-9901	Carte A Inc	157 356-9960	
37 296-4725	Carte A Inc	157 356-9960	
37 542-1521	Carte A Inc	157 356-9960	
37 364-5232	Carte A Inc	157 356-9960	
37 541-5649	Carte A Inc	157 356-9960	
37 739-2662	Carte A Inc	157 356-9960	
37 879-0030	Carte A Inc	157 356-9960	
37 541-3948	Carte A Inc	157 356-9960	
37 836-1511	Carte A Inc	157 356-9960	
37 569-4119	Carte A Inc	157 356-9960	
37 569-4782	Carte A Inc	157 356-9960	



\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**

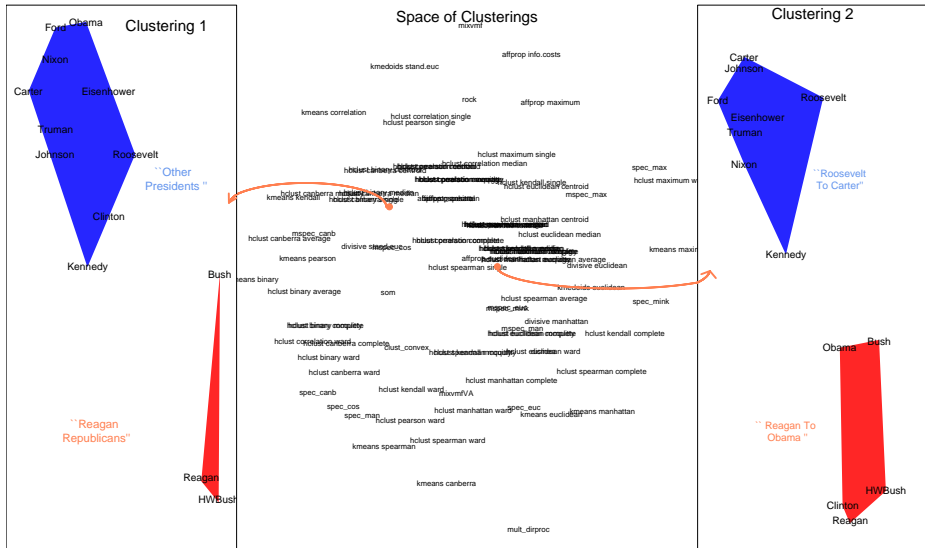
A New Strategy

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**
- 8 (Or, our new strategy: represent the entire bell space directly; no need to examine document contents)

Many Thousands of Clusterings, Sorted & Organized

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

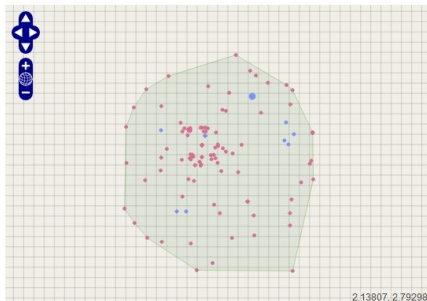


Software Screenshot

Size: 244 Files

Description: NSF - Updated Set

< > Number of Clusters ☒ 5 Clusters (Low) ☐ 15 Clusters (Medium) ☐ 30 Clusters (High) ☐ Discoverable



☒ Display History ☒ Display Method Points

Label	Coordinates	Clusters
an interesting clustering [Link]	-0.30819, 0.46229	5
methods-oriented clustering [Link]	0.84753, 1.42538	5

(*) Discoverable

Coordinates: 0.84753, 1.42538

Clusters: 5

Label [\[+\]](#) methods-oriented clustering

29.51% [\[Link\]](#)
research community health science public practice global political national urban
72
Label [\[+\]](#)

[View Detail](#)

27.46% [\[Link\]](#)
data economic markets policy survey models financial use not risk
67
Label [\[+\]](#)

[View Detail](#)

21.72% [\[Link\]](#)
human social science systems behavioral networks brain spatial complex dynamics
53
Label [\[+\]](#)

[View Detail](#)

15.16% [\[Link\]](#)
education students school learning creative skills teaching cognitive college teachers
37
Label [\[+\]](#)

[View Detail](#)

6.15% [\[Link\]](#)
language linguistic speech data speakers computer semantic cultural variation
15
documentation
Label [\[+\]](#)

[View Detail](#)

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
 - \Rightarrow 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
 - \Rightarrow 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
- \Rightarrow 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
- \Rightarrow 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
- $\text{Quality} = \text{mean}(\text{within cluster}) - \text{mean}(\text{between clusters})$

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
- \Rightarrow 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

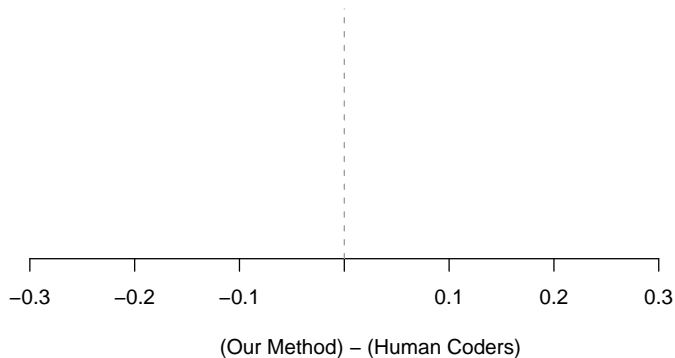
- What Are Humans Good For?

- They can't: keep many documents & clusters in their head
- They can: compare two documents at a time
- \Rightarrow Cluster quality evaluation: human judgement of document pairs

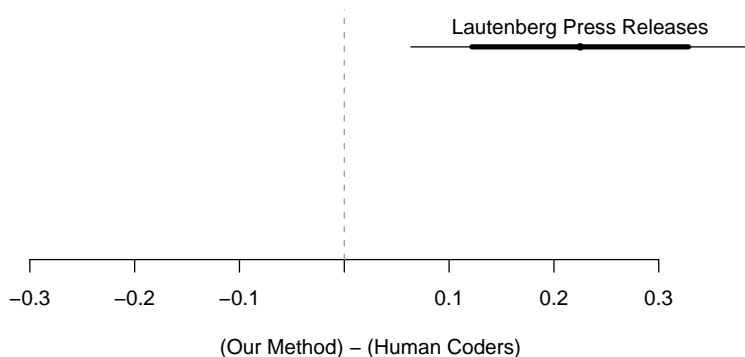
- Experimental Design to Assess Cluster Quality

- 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

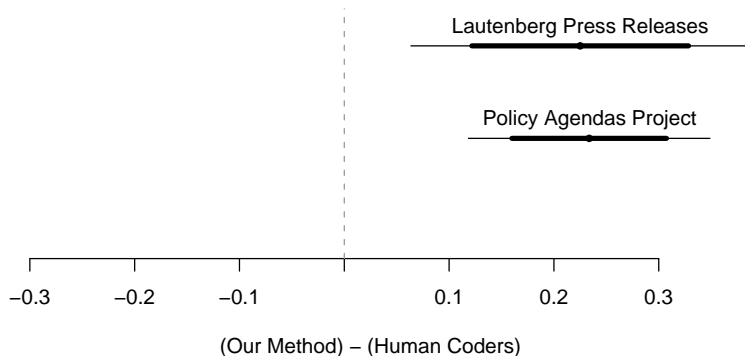


Evaluation 1: Cluster Quality



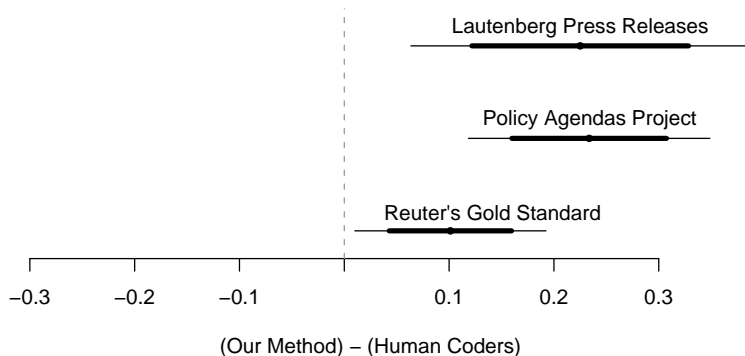
Lautenberg: 200 Senate Press Releases (appropriations, economy, education, tax, veterans, ...)

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

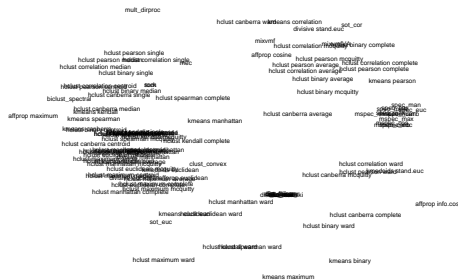
Evaluation 3: What Do Members of Congress Do?

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

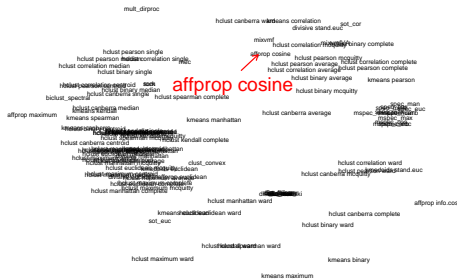
Evaluation 3: What Do Members of Congress Do?

- 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



Example Discovery



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

Example Discovery

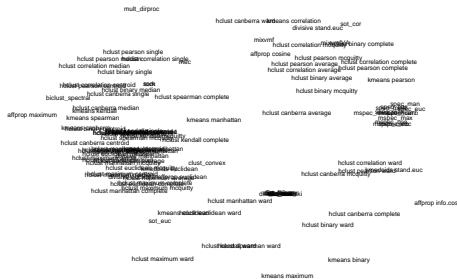


Red point: a **clustering** by
Affinity Propagation-Cosine
(Dueck and Frey 2007)

Close to:

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

Example Discovery



Space between methods:

Example Discovery



Space between methods:

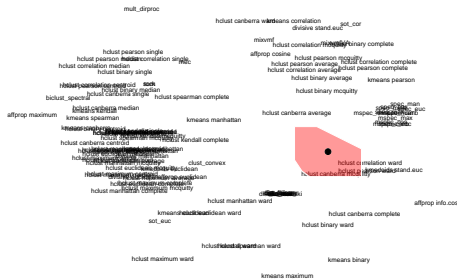
Example Discovery



Space between methods:
local cluster ensemble



Example Discovery



Mixture:

0.39 Hclust-Canberra-McQuitty

Example Discovery



Mixture:

0.39 Hclust-Canberra-McQuitty

0.30 Spectral clustering
Random Walk
(Metrics 1-6)

Example Discovery



Mixture:

0.39 Hclust-Canberra-McQuitty

0.30 Spectral clustering
Random Walk
(Metrics 1-6)

0.13 Hclust-Correlation-Ward

Example Discovery



Mixture:

0.39 Hclust-Canberra-McQuitty

0.30 Spectral clustering
Random Walk
(Metrics 1-6)

0.13 Hclust-Correlation-Ward

0.09 Hclust-Pearson-Ward

Example Discovery



Mixture:

0.39 Hclust-Canberra-McQuitty

0.30 Spectral clustering
Random Walk
(Metrics 1-6)

0.13 Hclust-Correlation-Ward

0.09 Hclust-Pearson-Ward

0.05 Kmediods-Cosine

Example Discovery



Mixture:

0.39 Hclust-Canberra-McQuitty

0.30 Spectral clustering
Random Walk
(Metrics 1-6)

0.13 Hclust-Correlation-Ward

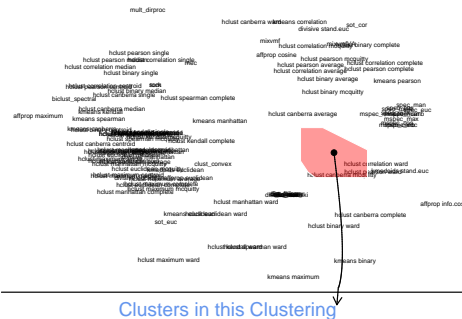
0.09 Hclust-Pearson-Ward

0.05 Kmediods-Cosine

0.04 Spectral clustering
Symmetric
(Metrics 1-6)

Clusters in this Clustering

Example Discovery



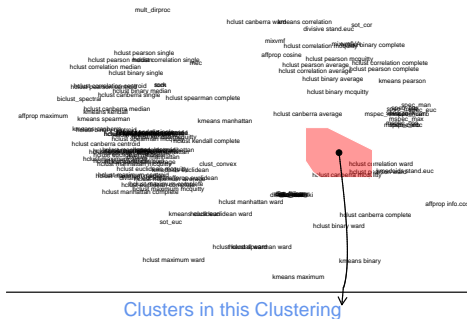
Credit Claiming Pork

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”

Mayhew

Example Discovery



Credit Claiming
Pork



Mayhew

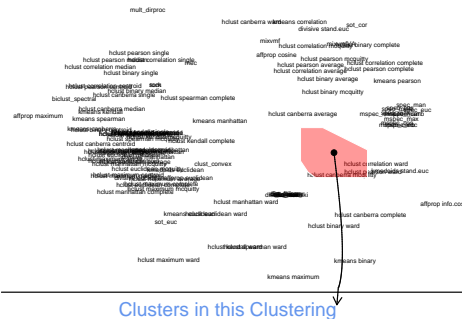
Credit Claiming
Legislation

Gary King (Harvard IQSS)

Credit Claiming, Legislation:

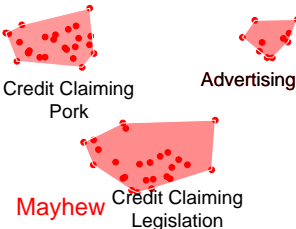
“As the Senate begins its recess, Senator Frank Lautenberg today pointed to a string of victories in Congress on his legislative agenda during this work period”

Example Discovery

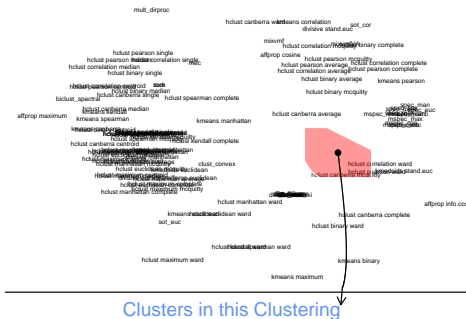


Advertising:

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

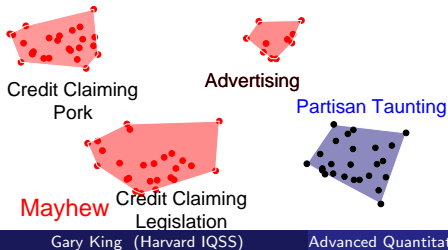


Example Discovery: Partisan Taunting

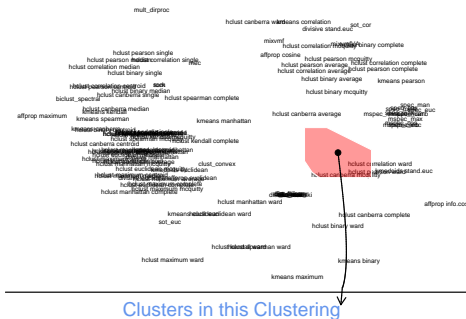


Partisan Taunting:

“Republicans Selling Out Nation on Chemical Plant Security”



Example Discovery: Partisan Taunting



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’”



Credit Claiming
Pork

Advertising

Partisan Taunting

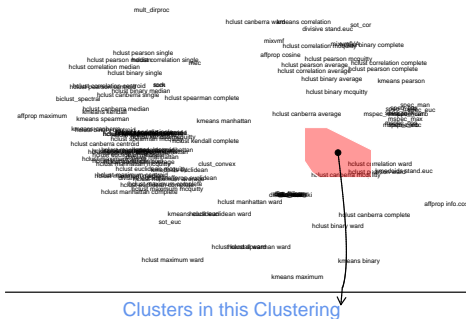
Mayhew

Credit Claiming Legislation

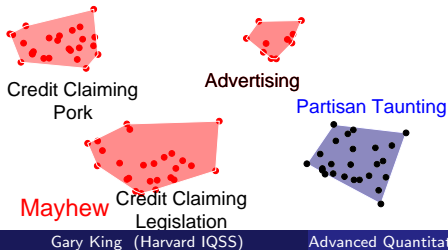
Gary King (Harvard IQSS)

Advanced Quantitative Research Methodology

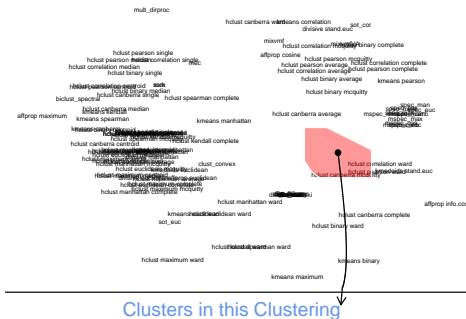
Example Discovery: Partisan Taunting



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



Example Discovery: Partisan Taunting



Credit Claiming
Pork

Advertising

Partisan Taunting

Mayhew

Credit Claiming
Legislation

Gary King (Harvard IQSS)

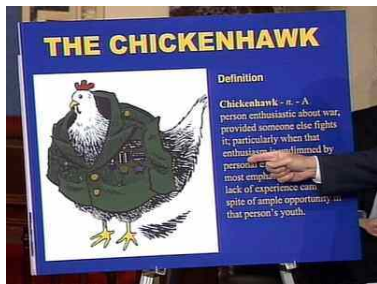
Advanced Quantitative Research Methodology

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

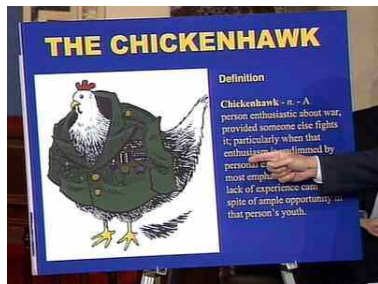


Sen. Lautenberg
on Senate Floor
4/29/04

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

In Sample Illustration of Partisan Taunting

Taunting ruins deliberation

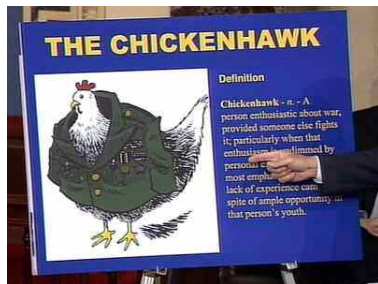


Sen. Lautenberg
on Senate Floor
4/29/04

- “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]

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

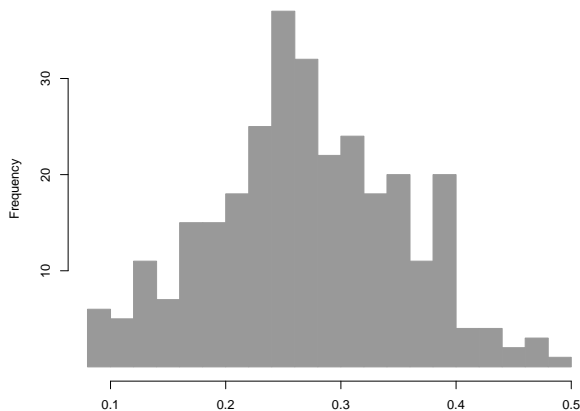
- Discovered using 200 press releases; 1 senator.
- Confirmed using 64,033 press releases; 301 senator-years.

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

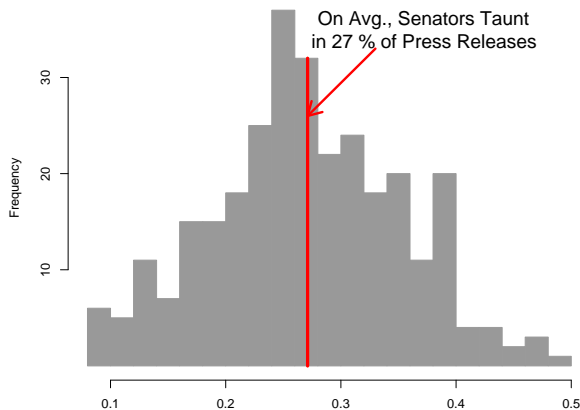
Out of Sample Confirmation of Partisan Taunting

- 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

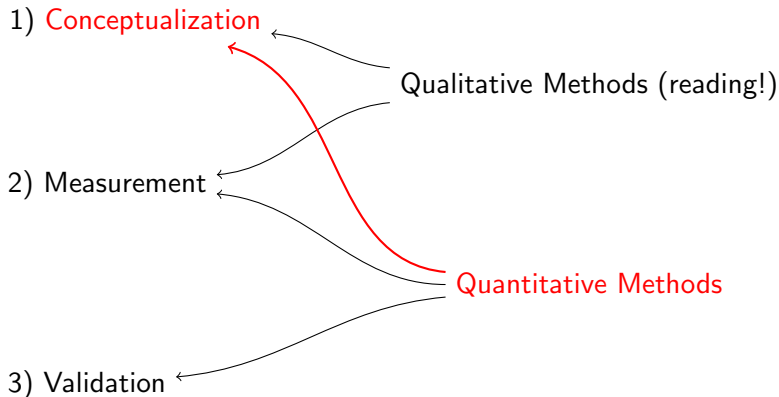


Out of Sample Confirmation of Partisan Taunting

- 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

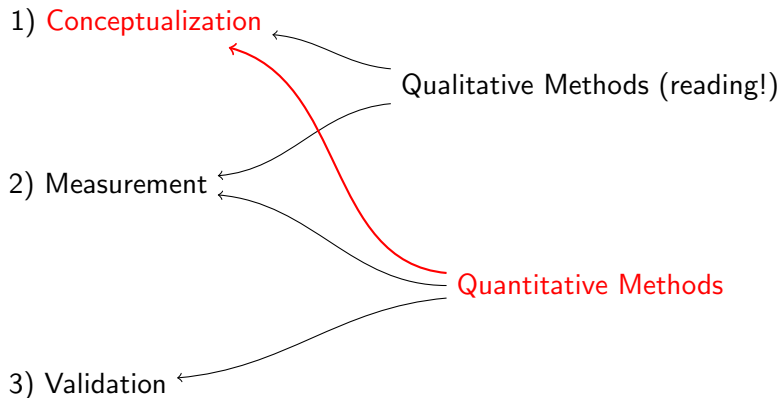


Quantitative Methods for Qualitative Conceptualization



Quantitative methods for conceptualization and discovery

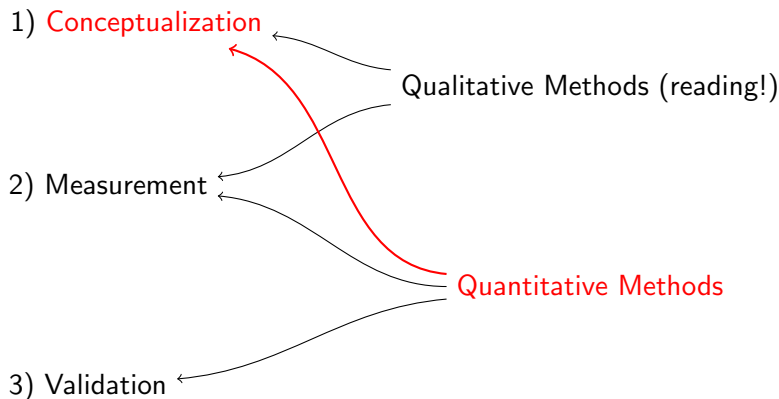
Quantitative Methods for Qualitative Conceptualization



Quantitative methods for conceptualization and discovery

- Few formal methods designed explicitly for conceptualization

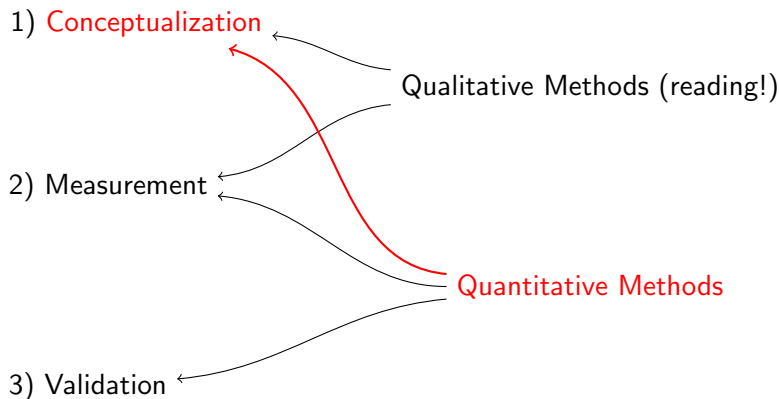
Quantitative Methods for Qualitative Conceptualization



Quantitative methods for conceptualization and discovery

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

Quantitative Methods for Qualitative Conceptualization



Quantitative methods for conceptualization and 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>