

Computer-Assisted Conceptualization

Gary King

Institute for Quantitative Social Science
Harvard University

Talk at Harvard Graduate School of Arts and Sciences, Alumni Day, 4/2/2011

¹Based on joint work with Justin Grimmer (Harvard ↔ 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).
- **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: Meaning Through Geography

Set of clusterings \approx

A list of unconnected addresses

wide at SuperPages.com

| | 195 | Car | C |
|--|---|---|---|
| Cartage New England Inc 28 Allen Ln Ipswich 01938..... 978 356-9960 | Carter F 34 Hibiscus Bldg 02133..... 617 327-1105 | Carter Nella E 323 Mainville Ave Box 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 Box 02186..... 617 698-5307 | |
| Cartagena Avish F Pleasant Box 02139..... 617 442-9780 | Francis S 134 Yankov W Ave 02132..... 617 323-6781 | Nick 21 Farwell Box 02114..... 617 267-5222 | |
| B Had 02134..... 617 361-5253 | Franklin & Anne 205 Mt Auburn Cam 02138..... 617 354-0798 | Nick & Debbi 196 Vermont Rd Newton 02459..... 617 527-0480 | |
| Jessica 50 Decatur Cha 02129..... 617 241-0152 | Fred 42 Hawthorn Hill 02136..... 617 524-3078 | Nicole..... 617 698-0713 | |
| Luzmila 124 Harvard Cam 02138..... 617 491-5621 | Fred 76 Howland Ave 02136..... 617 698-1343 | Norman G 38 Chickawhatch Dr 02125..... 617 822-1201 | |
| M 95 Howe Box 02132..... 617 323-9713 | G & B 8 Vardon Bldg 02134..... 617 434-8906 | P 40 Cranston Pl Box 02135..... 617 437-4754 | |
| Melvin 503 Green Cam 02139..... 617 576-1061 | G T 27 Franklin Ave Sun 02145..... 617 623-7121 | P E 501 E South S Box 02137..... 617 268-8213 | |
| Carte Nicholas 18 Appleton Boston 02114..... 617 695-6996 | Gayle 25 Franklin St 02133..... 617 823-0322 | P E 144 Hutchings Box 02131..... 617 427-9170 | |
| Cartagena D 4 Bradford Box 02133..... 617 338-0219 | George 225 Hudson Bldg 02134..... 617 367-9548 | P E 81 Boyden Ave 02138..... 617 968-8692 | |
| Carten Thos Jr Sr & Claire 17 Franklin Rd Mt 02136..... 617 698-6163 | Carter Hillside Assoc/Am 107 S Street Box 02111..... 617 456-1689 | Paul & Constance 114 Freeman Ave W Box 02131..... 617 325-2036 | |
| 17 445-5116 | Carter Harry F 26 Irving Ave Rt W Ave 02132..... 617 325-5465 | Paul E 501 E South S S Box 02137..... 617 268-4546 | |
| 17 822-2962 | Carter Hide Co Inc 167 Essex St 02131..... 617 542-7987 | Paul M 27 Crown Rd 02139..... 617 787-2115 | |
| 17 427-5712 | A Heber 617 442-5230 | Carter Pike Driving Inc 27 Beaver Ct Framingham 02702..... Wellesley Falls 781.235-0488 | |
| 17 569-2698 | Carter Hilary 41 Harvey Cam 02148..... 617 876-2750 | Carter Prudence 40 Franklin Waterman 02127..... 617 393-3782 | |
| 17 667-5190 | Horace 381 Walnut Ave Rosbury 02138..... 617 442-5307 | Prudence 40 Franklin Waterman 02127..... 617 926-7063 | |
| 17 569-1417 | Howard Jr 28 Neta Drive Box 02118..... 617 445-5532 | Roginald 106 Brookview Dorchester 02122..... 617 541-2843 | |
| 17 338-9110 | J Dan..... 617 354-2658 | Renee & Andrew 100 Walnut Box 02138..... 617 720-3765 | |
| 17 825-1953 | J 31 Chatham Box 02144..... 617 233-7990 | Carter Rice David 3450 Boston Publishing 163 Main Wilmington 01887 | |
| Carter Anne MD 1161 Beacon Bldg 02144..... 617 739-1022 | J 538 Harvard Box 02146..... 617 730-9483 | Ted Free-Dad 'I' & Thm..... 800 638-1671 | |
| Carter J M 1 Ipswich Pl Box 02146..... 617 735-8787 | J 775 The Pines West Rosbury 02132..... 617 323-5374 | Carl Eric Industrial Prod 613 Main Wilmington Ted Free-Dad 'I' & Thm..... 800 616-7447 | |
| 17 670-2078 | 1 Ipswich Pl Box 02146..... 617 735-8787 | Carl Free-Dad 'I' & Thm..... 800 648-7447 | |
| 17 621-9001 | 3410 Columbia Rd S Box 02137..... 617 464-1040 | Headquarters 613 Main Wilmington 02102 Carl..... 978 988-7447 | |
| 17 296-4725 | Carter J M 17 Ipswich Pl Box 02146..... 617 735-8787 | Ingalls Engine 163 Main Wilmington 01887 Carl..... 800 638-1673 | |
| 17 542-1521 | B E 18 Graduate Ave Mt 02136..... 617 296-6911 | Carter Richard 2079 Lawrence Ave Brighton 02215..... 617 987-0836 | |
| 17 364-5232 | Carter Barbara L MD Tufts-New England Medical Center Box 02111 Cam..... 617 436-0051 | Richard A 97 Mt Vernon Box 02106..... 617 566-7293 | |
| 17 541-5649 | Carter Becky Jo 02134..... 617 523-4368 | J 200 Conventry Pl Box 02136..... 617 267-0710 | |
| Carter Broadcasting Co 28 Park Pl Box 02134..... 617 423-0210 | Bernard J 122 Goodville E Box 02136..... 617 567-9430 | Carter Richard K M 13 Mather St Box 02137..... 617 268-0448 | |
| 17 739-2662 | Bibbiah 25 Midway Dr 02136..... 617 298-8713 | Roger 130 St Bourne Box 02131..... 617 424-6148 | |
| 17 879-0030 | Blair 26 Mt Vernon Box 02106..... 617 367-9031 | Roy 41 Concord Cam 02138..... 617 491-6115 | |
| 17 541-3948 | Carter & Baines Consultants Inc 28 Park Pl Box 02134..... 617 423-0210 | Royce 18 Sanyday Cha 02129..... 617 241-0418 | |
| 17 436-1511 | Carter C 200 Conventry Ave Box 02136..... 617 782-2118 | | |
| 17 569-4119 | C 218 Harvard Ave East Boston 02128..... 617 569-1545 | | |
| 800 569-8782 | C 109 Harvard Cam 02131..... 617 491-4822 | | |
| | C 109 Harvard Cam 02131..... 617 491-4822 | | |
| | C & M 43 Bernham Jan 02138..... 617 524-9558 | | |

Our Idea: Meaning Through Geography

Set of clusterings \approx

A list of unconnected addresses

wide at SuperPages.com

195

| Car | | Car | |
|-------------|---|-------------|--|
| 17 566-1282 | Cartage New England Inc 28 Allen Ln Essex 01828 | 17 327-1105 | Carter F. 514 Hickox Ave 02131 |
| 17 447-4101 | Cartagena Lydia 28 Sweet Briar 02131 | 17 437-7331 | Faye & Ricky 20 Columbia Ave 02131 |
| 17 257-9961 | Cartagena Avish F Beach Rd 02131 | 17 323-6781 | Francis S. 134 Temple W Ave 02131 |
| 17 566-1282 | B Had 02131 | 17 354-0798 | Franklin & Anne 705 Mt Auburn Cam 02131 |
| 17 364-5188 | Justica 50 Decatur Cha 02129 | 17 524-3078 | Fred 41 Haverhill Jam 02131 |
| 361-0380 | Luzella 124 Harvard Cam 02131 | 17 698-1343 | Fred 16 Haverhill Av 02131 |
| 17 566-4548 | M 90 Howe St 02131 | 17 434-8906 | G & B. 8 Vardon Ave 02131 |
| 17 628-8248 | Melvin 503 Green Cam 02129 | 17 623-7121 | G T 27 Frylands Av 02131 |
| 17 445-5116 | Carte Nicholas 18 Appleton Boston 02114 | 17 823-8322 | Gayle 25 Franklin St 02131 |
| 17 822-2962 | Cartier G. 4 Highland Ave 02131 | 17 522-3215 | Geo S 115 Mount Hill Nat 02131 |
| 17 427-5712 | Cartier H. 100 St 02131 | 17 367-9548 | George 125 Boston Ave 02131 |
| 17 569-2698 | A 202 Beulah Wy Haverhill 02131 | 17 456-1689 | Carter Holiday Assoc 107 S Street St 02111 |
| 17 667-5190 | A M 255 Main St 02131 | 17 325-5465 | Carter Harry F. 100 Burns Rd W Ave 02131 |
| 17 569-1417 | Adams 301 Carter St 02131 | 17 542-7987 | Carter Hide Co Inc 140 Burns Rd W Ave 02131 |
| 17 338-9110 | Allice 40 Market Cambridge 02139 | 17 876-2750 | Carter Hilary 41 Harvey Cam 02148 |
| 17 825-9195 | Andrew F. 42 Mt St 02131 | 17 442-5300 | Horace 301 Walnut Av Haverhill 02131 |
| 17 296-1293 | 1101 Beacon St 02144 | 17 442-5307 | Howard J. 28 New One Box 02118 |
| 17 670-2078 | 771 Newbury Boston 02116 | 17 445-5532 | J. Cam 41 Franklin Waterbury 02127 |
| 17 621-9001 | B. E. 10 Gladstone Ave 02131 | 17 232-7990 | Renée & Andrew 100 Broadway Center 02121 |
| 17 296-4725 | Cartier Barbara L. MD Tufts New England Medical Center Box 02111 | 17 730-9483 | Carter Rice Doan Building Design Publishing 163 Main Wilmington 01887 |
| 17 542-1521 | Cartier Becky 901 02131 | 17 323-5274 | Tom Free-Dad 'I' & Thom 600 638-1671 |
| 17 364-5232 | Bernard J. 301 Main St 02131 | 17 735-8787 | Tom Free-Dad 'I' & Thom 600 648-7447 |
| 17 541-5649 | Bernadette S. 20 Midway Ave 02131 | 17 492-1214 | Tom Free-Dad 'I' & Thom 600 648-7447 |
| 17 739-2662 | Cartier Broadcasting Co 50 Park Pl 02131 | 17 739-2193 | Tom Free-Dad 'I' & Thom 600 648-7447 |
| 17 879-0030 | Cartier C. 200 Gessert Av 02131 | 17 876-8841 | Tom Free-Dad 'I' & Thom 600 648-7447 |
| 17 541-3948 | Cartier C. 200 Gessert Av 02131 | 17 361-0773 | Tom Free-Dad 'I' & Thom 600 648-7447 |
| 17 436-1511 | C 210 Harvard Ave East Boston 02128 | 17 964-0435 | Tom Free-Dad 'I' & Thom 600 648-7447 |
| 17 569-4119 | C 109 Harvard Cam 02131 | 17 426-5094 | Tom Free-Dad 'I' & Thom 600 648-7447 |
| 17 809-8222 | C & M 41 Northgate Ave 02131 | 17 987-2163 | Tom Free-Dad 'I' & Thom 600 648-7447 |
| 17 809-8782 | C & M 41 Northgate Ave 02131 | 17 423-4334 | Tom Free-Dad 'I' & Thom 600 648-7447 |
| | | 17 282-1235 | Tom Free-Dad 'I' & Thom 600 648-7447 |
| | | 17 734-6109 | Tom Free-Dad 'I' & Thom 600 648-7447 |
| | | 17 265-8656 | Tom Free-Dad 'I' & Thom 600 648-7447 |
| | | 17 282-1593 | Tom Free-Dad 'I' & Thom 600 648-7447 |
| | | | Tom Free-Dad 'I' & Thom 600 648-7447 |



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 |
| 100 257-9961 | Cartagena Avish F Beach Rd 02139 | 617 442-9780 |
| 17 566-1282 | B Had 02134 | 617 361-5253 |
| 17 364-5188 | Lucille 174 Harvard Can 02139 | 617 491-5621 |
| 361-0380 | M 95 Howe Box 02135 | 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 | Carten Thos & Sr & Claire 1 Furlow St W 02135 | 617 338-9219 |
| 17 822-2962 | Carton S 10 Thompson Ln Mt 02136 | 617 696-6919 |
| 17 427-5712 | A Weber 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 255 Main St Av 02115 | 617 266-7153 |
| 17 569-1417 | Adams 301 Carter St Mt 02136 | 617 698-9074 |
| 17 338-1101 | Adams P 42 West St 02135 | 617 945-2711 |
| 17 825-1193 | Carte Anne MD 1101 Beacon St 02144 | 617 739-1022 |
| 17 296-1295 | Cartier Adhena 971 Newbury Boston 02116 | 617 536-6239 |
| 17 670-2078 | B C 10 Gladstone Av Mt 02136 | 617 296-6911 |
| 17 621-9001 | Cartier Barbara L MD Tufts New England Medical Center Box 02111 | 617 436-0951 |
| 17 296-4725 | Cartier Becky MD 02114 | 617 523-4368 |
| 17 542-1521 | Bernard J 301 Ashburne E Rd 02136 | 617 567-9430 |
| 17 364-5232 | Bibbith 25 Midway Dr 02134 | 617 298-8713 |
| 17 541-5649 | Billings 18 Waverley St 02136 | 617 367-9931 |
| 17 739-2662 | Cartier Broadcasting Co 50 Park Pl Box 02114 | 617 423-0210 |
| 17 879-0030 | Cartier C 2000 Cambridge St 02135 | 617 225-0200 |
| 17 541-3948 | C 210 Fremont Av East Boston 02128 | 617 782-2118 |
| 17 436-1511 | C 109 Harvard Can 02136 | 617 491-4822 |
| 17 569-4119 | C 8 M 41 Northgate Av 02134 | 617 524-9532 |
| 800 569-8782 | C & M 41 Northgate Av 02134 | 617 524-9532 |
| 17 327-1105 | Carter F 514 Hubbs Box 02131 | 617 327-1105 |
| 17 437-7331 | Faye & Ricky 20 Columbia Av Box 02136 | 617 437-7331 |
| 17 323-6781 | Francis S 134 Temple W Av 02132 | 617 323-6781 |
| 617 354-0798 | Franklin & Anne 705 Mt Auburn Can 02138 | 617 354-0798 |
| 617 524-3078 | Fred 41 Howard Av 02136 | 617 524-3078 |
| 617 698-1343 | Fred 16 Howley Av Mt 02136 | 617 698-1343 |
| 617 436-8906 | G & B 8 Vardon Box 02134 | 617 436-8906 |
| 617 623-7121 | G T 27 Fossil Av Mt 02135 | 617 623-7121 |
| 617 825-8322 | Gayle 25 Franklin St 02134 | 617 825-8322 |
| 617 522-3215 | Geo S 115 Main Mt Av 02136 | 617 522-3215 |
| 617 367-9548 | George 25 Hudson Box 02114 | 617 367-9548 |
| 617 456-1689 | Carter Hillside Assoc 107 S Street Box 02111 | 617 456-1689 |
| 617 325-5465 | Carter Harry F 100 Bayne Rd W Av 02132 | 617 325-5465 |
| 617 542-7987 | Carter Hide Co Inc 140 Boston St W 02132 | 617 542-7987 |
| 617 876-2750 | Carter Hilary 41 Harvey Can 02148 | 617 876-2750 |
| 617 442-5307 | Horace 301 Walnut Av Haverhill 02119 | 617 442-5307 |
| 617 445-5552 | Howard Jr 28 New One Box 02118 | 617 445-5552 |
| 617 354-2658 | J Can 15 Chatham St 02144 | 617 232-7990 |
| 617 730-9483 | J 538 Harvard Box 02144 | 617 730-9483 |
| 617 323-5274 | J 775 The Pines West Haverhill 02132 | 617 323-5274 |
| 617 735-8787 | J 1 Brooklyn Pl Box 02144 | 617 735-8787 |
| 617 464-1040 | 3410 Columbia Rd S Box 02137 | 617 464-1040 |
| 617 436-5353 | Carter J M Ornamental Ironworks 100 Franklin Falls 02174 | 617 436-5353 |
| 617 442-1775 | Carter J Veal Co 40 Newmarket St 02138 | 617 442-1775 |
| 617 492-1214 | 1573 Cambridge St Can 02136 | 617 492-1214 |
| 617 739-2193 | James 622 Foster Av Haverhill 02119 | 617 739-2193 |
| 617 876-8841 | J 101 East Star Rd Cambridge 02141 | 617 876-8841 |
| 617 361-0773 | J 34 Howley Rd Mt 02136 | 617 361-0773 |
| 617 964-0435 | Jane 14 Adams Rd Newton 02458 | 617 964-0435 |
| 617 426-9094 | John 1100 Waverley St 02136 | 617 426-9094 |
| 617 987-2163 | John 11 Mansfield St 02134 | 617 987-2163 |
| 617 423-4334 | John 207 Summer St 02135 | 617 423-4334 |
| 617 282-1235 | John 40 Howard St 02136 | 617 282-1235 |
| 617 734-6109 | James O 129 A Summit Av Box 02131 | 617 734-6109 |
| 617 265-8656 | J 29 Inverness Dr 02134 | 617 265-8656 |
| 617 282-1593 | K 17 Exposed Dr 02132 | 617 282-1593 |
| 617 267-6483 | 323 Main St Av Box 02115 | 617 267-6483 |
| 617 698-5307 | Nicholas S F 115 Randolph Av Mt 02136 | 617 698-5307 |
| 617 267-5222 | Nick 21 Furlow Box 02114 | 617 267-5222 |
| 617 527-0480 | 136 Hermit Rd Newton 02459 | 617 527-0480 |
| 617 698-0713 | Norman G 38 Chickadee Dr 02125 | 617 698-0713 |
| 617 822-1201 | 38 Chickadee Dr Box 02125 | 617 822-1201 |
| 617 427-4754 | P 44 Haverhill Box 02135 | 617 427-4754 |
| 617 268-8213 | P E 501 E South S Box 02137 | 617 268-8213 |
| 617 427-9170 | P L 44 Haverhill Box 02135 | 617 427-9170 |
| 617 968-8692 | P R 91 Bayne Av 02138 | 617 968-8692 |
| 617 325-3034 | 114 Adams Av W Mt 02132 | 617 325-3034 |
| 617 268-4546 | Paul F 501 E South S Box 02137 | 617 268-4546 |
| 617 787-2115 | Paul M 27 Union St 02139 | 617 787-2115 |
| 617 235-0488 | Carter Pike Driving Inc 27 Beaver Ct Framingham 02702 | 617 235-0488 |
| 617 393-3782 | Carter Prudence 40 Franklin Waterbury 02172 | 617 393-3782 |
| 617 926-7063 | Prudence 40 Franklin Waterbury 02172 | 617 926-7063 |
| 617 541-2843 | Reginald 100 Broadview Dr 02122 | 617 541-2843 |
| 617 720-3765 | Carter R 10 Walnut Box 02138 | 617 720-3765 |
| 800 638-1671 | Carter Rice David Building Division Publishing 163 Main Wilmington 01887 | 800 638-1671 |
| 800 616-7447 | Carter Richard A 100 Broadview Dr 02122 | 800 616-7447 |
| 800 648-7447 | Carter Richard A 100 Broadview Dr 02122 | 800 648-7447 |
| 978 988-7447 | Carter Richard A 100 Broadview Dr 02122 | 978 988-7447 |
| 800 638-1673 | Carter Richard A 100 Broadview Dr 02122 | 800 638-1673 |
| 617 987-0836 | Carter Richard A 2075 Carleton Av Brighton 02111 | 617 987-0836 |
| 617 566-7293 | Carter Richard A 100 Broadview Dr 02122 | 617 566-7293 |
| 617 267-0710 | Carter Richard A 100 Broadview Dr 02122 | 617 267-0710 |
| 617 268-0468 | Carter Richard K 123 Merwin St Box 02137 | 617 268-0468 |
| 617 864-1535 | Robert L 175 Rockwood Av Can 02141 | 617 864-1535 |
| 617 424-6148 | Royce 130 St Brnagh Box 02131 | 617 424-6148 |
| 617 491-6115 | Royce 130 St Brnagh Box 02131 | 617 491-6115 |
| 617 241-9418 | Royce 18 Sandway Cir 02129 | 617 241-9418 |



\approx We develop a (conceptual) geography of clusterings

A New Strategy

Make it easy to choose best clustering from millions of choices

A New Strategy

Make it easy to choose best clustering from millions of choices

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

A New Strategy

Make it easy to choose best clustering from millions of choices

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

A New Strategy

Make it easy to choose best clustering from millions of choices

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

A New Strategy

Make it easy to choose best clustering from millions of choices

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

A New Strategy

Make it easy to choose best clustering from millions of choices

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

A New Strategy

Make it easy to choose best clustering from millions of choices

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

A New Strategy

Make it easy to choose best clustering from millions of choices

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

A New Strategy

Make it easy to choose best clustering from millions of choices

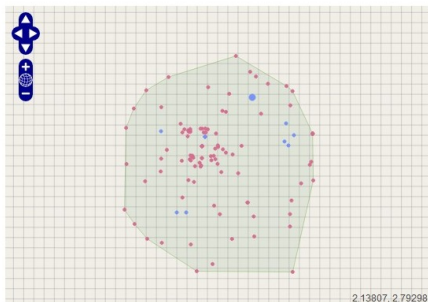
- 1 **Code text as numbers** (in one *or more* of several ways)
- 2 **Apply all clustering methods we can find** to the data — each representing different (unstated) substantive assumptions (<15 mins)
- 3 (Too much for a person to understand, but organization will help)
- 4 Develop an **application-independent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

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%
72 research community health science public practice global political national urban
Label [+]

27.46%
67 data economic markets policy survey models financial use not risk
Label [+]

21.72%
53 human social science systems behavioral networks brain spatial complex dynamics
Label [+]

15.16%
37 education students school learning creative skills teaching cognitive college teachers
Label [+]

6.15%
15 language linguistic speech data speakers computer semantic cultural variation
documentation
Label [+]

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)

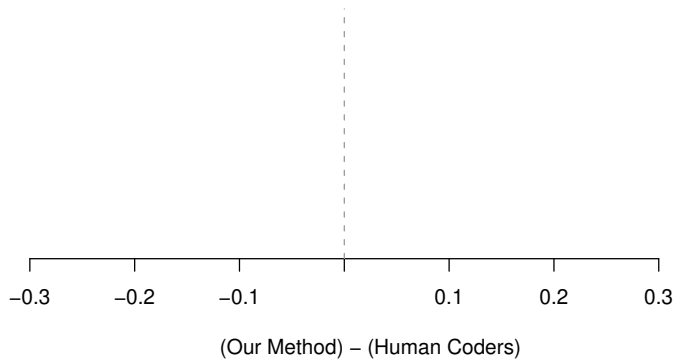
Evaluation 1: Cluster Quality

- **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**

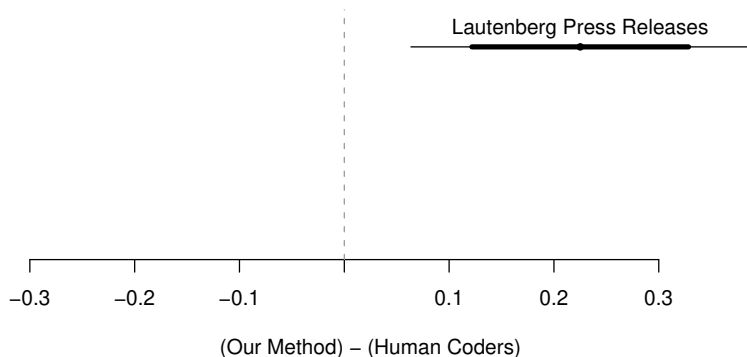
Evaluation 1: Cluster Quality

- **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

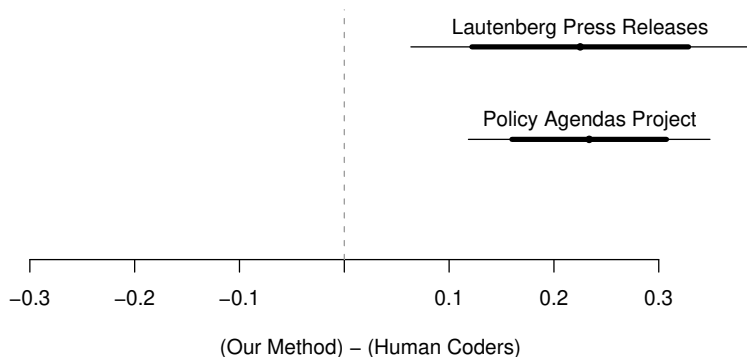


Evaluation 1: Cluster Quality



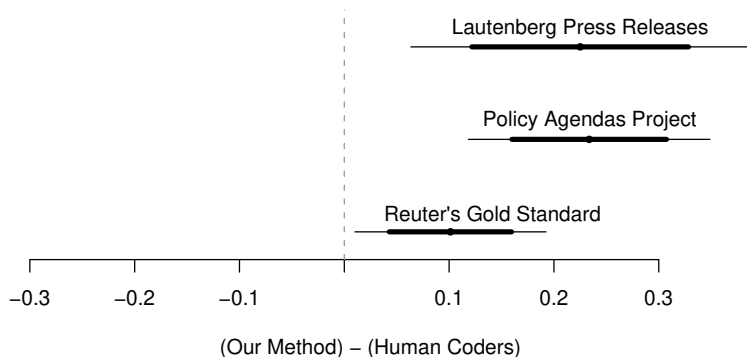
Lautenberg: 200 Senate Press Releases (appropriations, economy, education, tax, veterans, ...)

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



Space between methods:
local cluster ensemble

Example Discovery



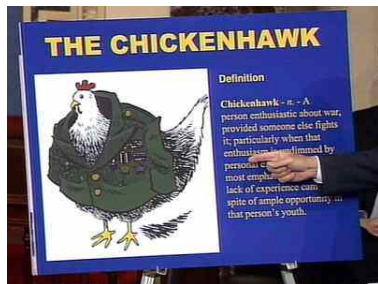
Mixture:

0.39 Hclust-Canberra-McQuitty

0.30 Spectral clustering
Random Walk
(Metrics 1-6)

0.13 Hclust-Correlation-Ward

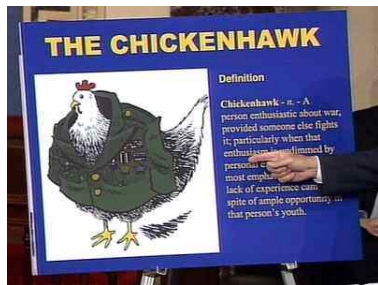
Taunting ruins deliberation



Sen. Lautenberg
on Senate Floor
4/29/04

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

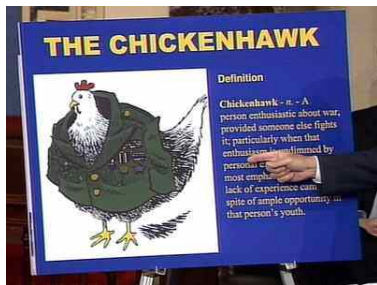
Taunting ruins deliberation



Sen. Lautenberg
on Senate Floor
4/29/04

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

Taunting ruins deliberation



Sen. Lautenberg
on Senate Floor
4/29/04

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

Out of Sample Confirmation of Partisan Taunting

- Discovered using 200 press releases; 1 senator.

Out of Sample Confirmation of Partisan Taunting

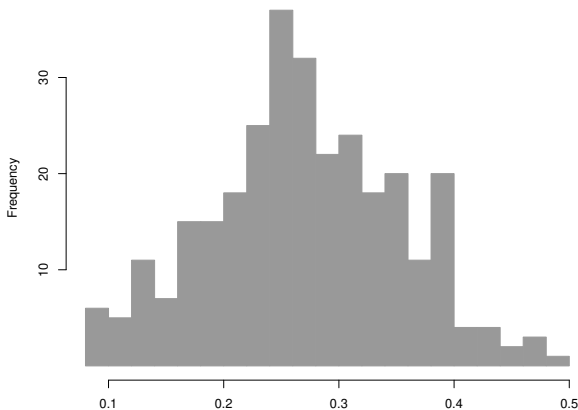
- Discovered using 200 press releases; 1 senator.
- Confirmed using 64,033 press releases; 301 senator-years.

Out of Sample Confirmation of Partisan Taunting

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

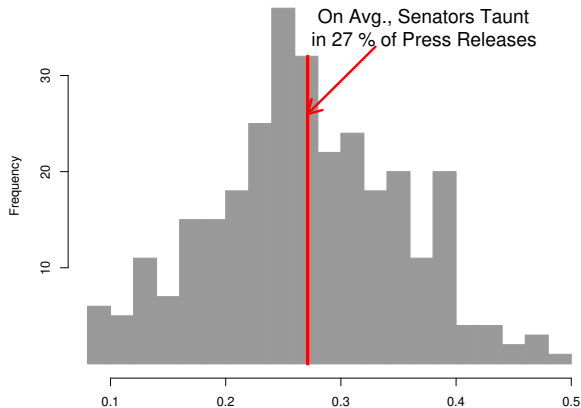
Out of Sample Confirmation of Partisan Taunting

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

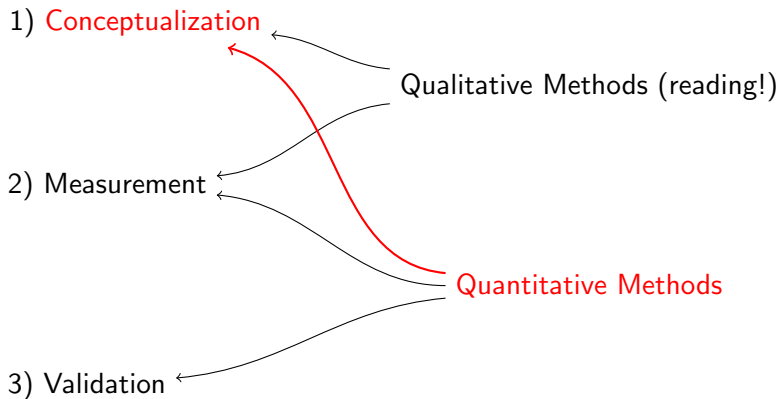


Out of Sample Confirmation of Partisan Taunting

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

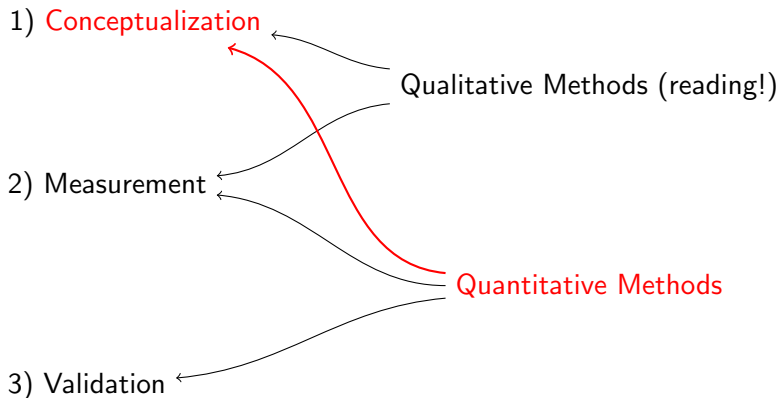


Quantitative Methods for Qualitative Conceptualization



Computer-Assisted Methods for conceptualization and discovery

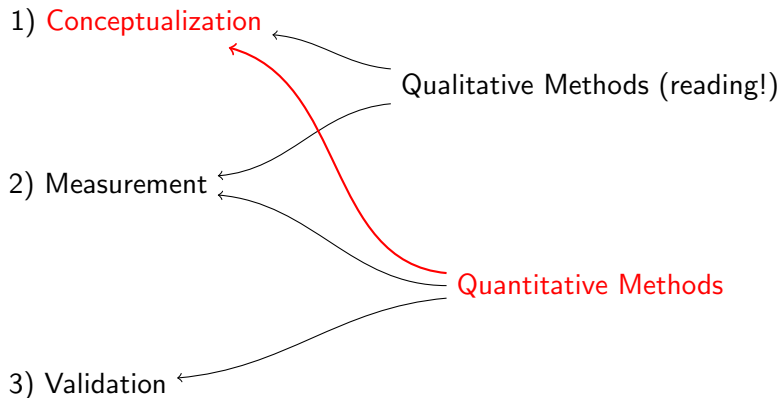
Quantitative Methods for Qualitative Conceptualization



Computer-Assisted Methods for conceptualization and discovery

- Methods designed explicitly for conceptualization

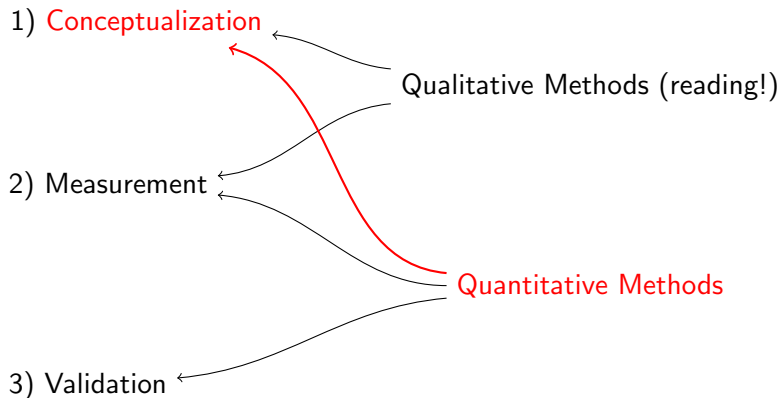
Quantitative Methods for Qualitative Conceptualization



Computer-Assisted Methods for conceptualization and discovery

- Methods designed explicitly for conceptualization
- The end of quantitative v qualitative debates

Quantitative Methods for Qualitative Conceptualization



Computer-Assisted Methods for conceptualization and discovery

- Methods designed explicitly for conceptualization
- The end of quantitative v qualitative debates
- Evaluation methods measure progress in discovery

For more information



<http://GKing.Harvard.edu>