

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

Gary King

Institute for Quantitative Social Science  
Harvard University

Talk at the Triangle Political Methods Group, Duke University on 2/11/2010

Joint work with Justin Grimmer (Harvard ↔ Stanford)

# Some context for related technology

- <http://ow.ly/14hDU> (play after ad)
- <http://ow.ly/14h36> (play 10:00-12:06)

# A Method for Conceptualization

- Systematic method for computer-assisted conceptualization from text

# A Method for Conceptualization

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

# A Method for Conceptualization

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

# A Method for Conceptualization

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

# Why Johnny Can't Classify (Optimally)

# Why Johnny Can't Classify (Optimally)

- Clustering seems easy; its not!



# Why Johnny Can't Classify (Optimally)

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

# Why Johnny Can't Classify (Optimally)

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

# Why Johnny Can't Classify (Optimally)

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

# Why Johnny Can't Classify (Optimally)

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

# Why Johnny Can't Classify (Optimally)

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

# Why Johnny Can't Classify (Optimally)

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

# Why Johnny Can't Classify (Optimally)

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

# Why Johnny Can't Classify (Optimally)

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



# Why HAL Can't Classify Either

# Why HAL Can't Classify Either

- Large quantitative literature on **cluster analysis**

# Why HAL Can't Classify Either

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

# Why HAL Can't Classify Either

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

# Why HAL Can't Classify Either

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

# Why HAL Can't Classify Either

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

# Why HAL Can't Classify Either

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

# Why HAL Can't Classify Either

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



# Why HAL Can't Classify Either

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

# Why HAL Can't Classify Either

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

# Why HAL Can't Classify Either

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

# If Ex Ante doesn't work, try Ex Post

# If Ex Ante doesn't work, try Ex Post

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

# If Ex Ante doesn't work, try Ex Post

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

# If Ex Ante doesn't work, try Ex Post

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

# If Ex Ante doesn't work, try Ex Post

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



# If Ex Ante doesn't work, try Ex Post

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

# If Ex Ante doesn't work, try Ex Post

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

# If Ex Ante doesn't work, try Ex Post

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

# Our Idea: Meaning Through Geography

Set of clusterings

# Our Idea: Meaning Through Geography

Set of clusterings  $\approx$

A list of unconnected addresses

wide at [SuperPages.com](http://SuperPages.com)

	195	Car	C
<b>Cartage New England Inc</b> 28 Allen Ln Ipswich 01938..... 978 356-9960	<b>Carter F</b> 34 Hibiscus Bldg 02133..... 617 327-1105	<b>Carter Nella E</b> 323 Main St 02115..... 617 267-6483	
<b>Cartagena Lydia</b> 20 Sweet Box 02331..... 617 323-7639	<b>Faye &amp; Ricky</b> 20 Columbia Ave Box 02136..... 617 437-7331	<b>Nicholas S F</b> 115 Randolph Ave Mill 02186..... 617 698-5307	
<b>Cartagena Avish</b> F Pleasant Box 02139..... 617 442-9780	<b>Francis S</b> 134 Yankov W Ave 02132..... 617 323-6781	<b>Nick 21 Farwell Box 02114..... 617 267-5222</b>	
<b>B Hed 02134</b> ..... 617 361-5253	<b>Franklin &amp; Anne</b> 201 Mt Auburn Cam 02138..... 617 354-0798	<b>Nick &amp; Debbi</b> 196 Herold Rd Newton 02459..... 617 527-0480	
<b>Jessica</b> 50 Decatur Cha 02129..... 617 241-0152	<b>Fred 42 Howland Elm 02136..... 617 524-3078</b>	<b>Nicole</b> ..... 617 698-0713	
<b>Luzmila</b> 124 Harvard Cam 02136..... 617 491-5621	<b>Fred 16 Howland Ave Mill 02136..... 617 698-1343</b>	<b>Norman G</b> 38 Chickawhatch Dr 02125..... 617 822-1201	
<b>M 90 Howe Box 02132</b> ..... 617 323-9713	<b>G &amp; B</b> 8 Vardon Dr 02134..... 617 434-8906	<b>P 40 Cranford Pl Box 02135</b> ..... 617 437-4754	
<b>Melvin</b> 503 Green Cam 02139..... 617 576-1061	<b>G T 27 Franklin Ave Sun 02145..... 617 623-7121</b>	<b>P E 501 E South S Box 02137</b> ..... 617 268-8213	
<b>Carl Nicholas</b> 18 Appleton Boston 02114..... 617 695-6996	<b>George</b> 125 Hudson Box 02134..... 617 367-9548	<b>P E 14 Hutchings Box 02131</b> ..... 617 427-9170	
<b>Carlton</b> 0 4 Bradford Box 02133..... 617 338-9219	<b>Carter Hillside Assoc</b> 107 S Street Box 02111..... 617 456-1689	<b>Paul &amp; Constance</b> 114 Franklin St W Box 02131..... 617 325-2036	
<b>Carten Thos J Sr &amp; Claire</b> 1 Franklin St Mill 02136..... 617 698-6163	<b>Carter Harry F</b> 30 Bayview Rd W Box 02132..... 617 325-5465	<b>Paul M 27 Crown St 02139</b> ..... 617 787-2115	
<b>17 445-5116</b> Thomas & Kathleen 50 Thompson Ln Mill 02136..... 617 696-6919	<b>Carter Hide Co Inc</b> 26 Burwell St W Box 02132..... 617 542-7987	<b>Paul M 27 Crown St 02139</b> ..... 617 787-2115	
<b>17 822-2962</b> Carter A Box 02133..... 617 229-2257	<b>Carter Hilary 41 Harvey Cam 02148</b> ..... 617 876-2750	<b>Prudence</b> 40 Franklin Waterlton 02127..... 617 393-3782	
<b>17 427-5712</b> A Heber..... 617 442-5230	<b>Horace</b> 301 Walnut St Rosbury 02139..... 617 442-5307	<b>Prudence</b> 40 Franklin Waterlton 02127..... 617 926-7063	
<b>17 569-2698</b> A 201 Beulah Wy Rosbury 02139..... 617 442-1219	<b>Howard Jr</b> 28 Nona Drive Box 02118..... 617 445-5532	<b>Roginald</b> 106 Brookview Dorchester 02122..... 617 541-2843	
<b>17 667-5190</b> A 201 Beulah Wy Rosbury 02139..... 617 442-1219	<b>J Dan</b> ..... 617 354-2658	<b>Renee &amp; Andrew</b> 10 Walnut Box 02118..... 617 720-3765	
<b>17 569-1412</b> A 201 Beulah Wy Rosbury 02139..... 617 442-1219	<b>J 21 Chatham Box 02144</b> ..... 617 232-7990	<b>Rice Doreen Publishing</b> 163 Main Wilmington 01887 Tel Free-Dial '9 & Then..... 800 638-1671	
<b>17 338-9110</b> Alicia 40 Market Cambridge 02139..... 617 945-2711	<b>J 538 Harvard Box 02146</b> ..... 617 730-9483	<b>Tel Free-Dial '9 &amp; Then</b> ..... 800 616-7447	
<b>17 825-1953</b> Andrew F 42 West St Box 02133..... 617 625-7623	<b>J 775 The Pines West Rosbury 02132</b> ..... 617 323-5374	<b>Tel Free-Dial '9 &amp; Then</b> ..... 800 648-7447	
<b>17 296-1593</b> 1101 Beacon Box 02144..... 617 739-1022	<b>J Brookline Pl Box 02146</b> ..... 617 735-8787	<b>Headquarters</b> 611 Main Wilmington 01887 Tel Free-Dial '9 & Then..... 800 648-1673	
<b>17 670-2078</b> B E 18 Graduate Ave Mill 02136..... 617 296-6911	<b>Carter J M</b> 3410 Columbia Rd S Box 02137..... 617 464-1040	<b>Call</b> ..... 800 616-7447	
<b>17 621-9001</b> Carter Barbara L MD Tufts-New England Medical Center Box 02111 Cam..... 617 436-0051	<b>Carter J M Ornamental Ironworks</b> Pondside Falls 02136..... 617 876-5353	<b>Call</b> ..... 800 648-7447	
<b>17 296-4725</b> Carter Becky Jo 02134..... 617 523-4368	<b>Carter J Veal Co</b> 40 Howland Elm 02136..... 617 442-1775	<b>Call</b> ..... 978 988-7447	
<b>17 542-1521</b> Bernard J 132 Goodhue E Box 02136..... 617 567-9430	<b>Carter James</b> 157 Cambridge St Cam 02136..... 617 492-1214	<b>Call</b> ..... 800 648-1673	
<b>17 364-5232</b> Bibbiah 25 Midway Dr 02134..... 617 298-8713	<b>James</b> 402 Foster St Rosbury 02136..... 617 739-2193	<b>Carter Richard</b> 2079 Lowellville Ave Brighton 02111..... 617 982-0836	
<b>17 541-5249</b> Bibbiah 25 Midway Dr 02134..... 617 298-8713	<b>James L</b> 34 Rosbury Rd Mill 02134..... 617 361-0773	<b>Carter Richard A MD</b> 130 Canterbury St Box 02136..... 617 267-0710	
<b>17 739-2662</b> Carter Broadcasting Co 50 Park Pl Box 02134..... 617 423-0210	<b>Janice</b> 14 Adams Rd Newton 02459..... 617 564-0435	<b>Carter Richard K</b> 23 Mather S Box 02137..... 617 268-0448	
<b>17 879-0030</b> Carter C 200 Casswell Ave Box 02135..... 617 782-2118	<b>Jeffrey</b> 41 Warren St Box 02134..... 617 424-5094	<b>Roy 130 St Bourne Box 02111</b> ..... 617 424-6148	
<b>17 436-1511</b> C 210 Harvard Ave East Boston 02128..... 617 569-1545	<b>John 111 Mansfield Pl 02134</b> ..... 617 987-2163	<b>Roy 41 Concord Cam 02138</b> ..... 617 491-6115	
<b>17 569-4119</b> C 109 Harvard Cam 02136..... 617 491-4822	<b>John 107 Summer Box 02129</b> ..... 617 423-4334	<b>Royce</b> 18 Safford Cha 02129..... 617 241-0418	
<b>800 569-4782</b> C 109 Harvard Cam 02136..... 617 491-4822	<b>John 40 Howland Elm 02136</b> ..... 617 262-1235		
	<b>John D 129 A Summit Ave Box 02133</b> ..... 617 734-6109		
	<b>J 29 Howland Elm 02136</b> ..... 617 265-8456		
	<b>K 17 Concord Dorchester 02122</b> ..... 617 282-1593		

# Our Idea: Meaning Through Geography

Set of clusterings  $\approx$

A list of unconnected addresses

wide at [SuperPages.com](http://SuperPages.com)

195

Car

C

17 566-1282	Cartage New England Inc 28 Allen St Ipswich 01938	978 356-9960	Carter F. 514 Hickox Ave 02131	617 327-1105	Carter Nella E 323 Marchant Ave Box 02115	617 267-6483
17 447-4101	Cartagena Lydia 28 Sweet Briar 02131	617 323-7639	Faye & Ricky 20 Columbia Ave Box 02136	617 437-7331	Nicholas S F 115 Randolph Ave Box 01386	617 698-5307
100 257-9961	Cartagena Avish F Beach Rd 02139	617 442-9780	Francis S. 134 Temple W Ave 02132	617 323-6781	Nick & Debbi 215 Fyfield Ave 02116	617 267-5222
17 566-1282	B Had 02136	617 361-5253	Franklin & Anne 705 Mt Auburn Cam 02138	617 354-0798	Norman G 196 Hermit Rd Newton 02459	617 527-0480
17 364-5188	Justicia 50 Decatur Cha 02129	617 241-0152	Fred 41 Haverhill Aven 02136	617 524-3078	Nick & Debbi 38 Chickadee Rd 02125	617 822-1203
361-0380	Luzmila 124 Harvard Cam 02136	617 491-5621	Fred W. 96 Valley St 02136	617 698-1343	P E 501 E South St Box 02137	617 268-8213
17 566-4548	M 95 Howe Box 02132	617 323-9713	G & B. 8 Vardon Ave 02134	617 436-8906	P L 44 Hutchings Box 02131	617 427-9170
17 628-8248	Melvin 503 Green Cam 02139	617 576-1061	Gayle 25 Franklin St 02134	617 823-8322	P R 91 Brewer Ave 02138	617 968-8692
17 445-5116	Carte Nicholas 18 Appleton Boston 02114	617 695-6996	Geo S 115 Mass Hill Rd Box 02138	617 522-3215	Paul & Constance 114 Adams Ave W Mass 02133	617 325-3034
17 822-2962	Cartagena O 4 Bradford Box 02133	617 338-9219	George 120 Nones St 02134	617 367-9548	Paul M 501 E South St Box 02137	617 268-4546
17 427-5712	Carten Thos J Sr & Claire 1 Fyfield St Mt 02116	617 698-6163	Carter Holiday Assoc 107 S Street Box 02111	617 456-1689	Paul M 27 Union St 02135	617 787-2115
17 569-2698	Carte Thos & Kathleen 50 Thompson Ln Mt 02136	617 696-6919	Carter Harry F 30 Burns Rd Rt W Ave 02132	617 325-5465	Prudence 40 Franklin Waterfront 02172	617 393-3782
17 667-5190	Carte A A 202 Pioneer Av Cambridge 02142	617 492-4174	Carter Hide Co Inc 161 Elm St 02148	617 542-7987	Reginald 100 Brookview Center 02123	617 541-2843
17 569-1417	Adams 361 Carter St Mt 02136	617 698-9074	Carter Hilary 41 Harvey Cam 02148	617 876-2750	Renee & Andrew 41 Main Wilmington 01857	800 638-1673
17 338-9110	Alice 40 Market Cambridge 02139	617 945-2711	Horace 301 Walnut Av Roxbury 02119	617 442-5307	Richard A 2077 Carver Ave Brighton 02115	617 982-0836
17 825-9195	Andrew F 42 Mt St 02135	617 625-7623	Howard Jr 28 New One Box 02118	617 445-5532	Richard A 47 Mt Vernon Box 02136	617 566-7293
17 296-1293	Carter Anne MD 1101 Beacon Bldg 02144	617 739-1022	J 58 Harvard St 02144	617 232-7990	Richard A 1200 Cambridge St 02136	617 967-0710
17 670-2078	Carter Arthur 971 Newbury Boston 02116	617 536-6229	J & Con 41 Chatham St 02144	617 232-7990	Richard A 1200 Cambridge St 02136	617 967-0710
17 621-9001	Carter Barbara L MD Tufts New England Medical Center Box 02111	617 296-6911	Jason 4775 The Pines West Roxbury 02132	617 730-9483	Richard A 1200 Cambridge St 02136	617 967-0710
17 296-4725	Carter Becky Joy 02134	617 636-0951	Jeffrey J 1300 Cambridge St Cam 02138	617 442-1775	Richard A 1200 Cambridge St 02136	617 967-0710
17 542-1521	Bernard J 301 Washington St 02136	617 523-4368	James 1503 Cambridge St Cam 02138	617 492-1214	Richard A 1200 Cambridge St 02136	617 967-0710
17 364-5232	Bibbath 25 Midway Rd 02136	617 298-8713	James 452 Foster Av Roxbury 02132	617 739-2193	Richard A 1200 Cambridge St 02136	617 967-0710
17 541-5649	Bliss 30 W Ames Ave 02136	617 367-9931	James 31 East Star Rd Cambridge 02141	617 876-8841	Richard A 1200 Cambridge St 02136	617 967-0710
17 739-2662	Carter Broadcasting Co 50 Park Pl Box 02136	617 423-0210	Jane L. 34 Rosbury Rd Mt 02136	617 361-0773	Richard A 1200 Cambridge St 02136	617 967-0710
17 879-0030	Carter C 73 East C Cam 02141	617 225-2020	Jane 14 Rosbury Rd Mt 02136	617 361-0773	Richard A 1200 Cambridge St 02136	617 967-0710
17 541-3948	Carter C 2000 Cambridge St 02136	617 782-2118	Jane 14 Rosbury Rd Mt 02136	617 361-0773	Richard A 1200 Cambridge St 02136	617 967-0710
17 436-1511	C 210 Harvard Av East Boston 02128	617 569-1545	John 11 Mansfield St 02134	617 987-2163	Richard A 1200 Cambridge St 02136	617 967-0710
17 569-6119	C 109 Harvard Cam 02136	617 491-4822	John 207 Summer St 02125	617 423-4334	Richard A 1200 Cambridge St 02136	617 967-0710
800 622-8782	C 8 & M 41 Northgate Ave 02134	617 524-9558	John 40 Westfield St 02129	617 282-1235	Roger 130 Stoughton St 02130	617 491-6115
			John 129 A Summit Av Box 02133	617 734-6109	Royce 18 Salisbury Cha 02129	617 241-9418
			K 775 The Pines West Roxbury 02132	617 765-8656		
			K 775 The Pines West Roxbury 02132	617 282-1593		





# A New Strategy

Make it easy to choose best clustering from millions of choices



# A New Strategy

Make it easy to choose best clustering from millions of choices

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

# A New Strategy

Make it easy to choose best clustering from millions of choices

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

# A New Strategy

Make it easy to choose best clustering from millions of choices

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

# A New Strategy

Make it easy to choose best clustering from millions of choices

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

# A New Strategy

Make it easy to choose best clustering from millions of choices

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

# A New Strategy

Make it easy to choose best clustering from millions of choices

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

# A New Strategy

Make it easy to choose best clustering from millions of choices

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





# Application-Independent Distance Metric: Axioms

# Application-Independent Distance Metric: Axioms

- Metric based on 3 assumptions

# Application-Independent Distance Metric: Axioms

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

# Application-Independent Distance Metric: Axioms

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

# Application-Independent Distance Metric: Axioms

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

# Application-Independent Distance Metric: Axioms

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

# Application-Independent Distance Metric: Axioms

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

# Evaluating the Performance of Our Method



# Evaluating the Performance of Our Method

- Goals:

# Evaluating the Performance of Our Method

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

# Evaluating the Performance of Our Method

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

# Evaluating the Performance of Our Method

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

# Evaluating the Performance of Our Method

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

# Evaluating the Performance of Our Method

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

# Evaluating the Performance of Our Method

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

# Evaluating the Performance of Our Method

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



# Evaluation 1: Cluster Quality

# Evaluation 1: Cluster Quality

- What Are Humans Good For?

# Evaluation 1: Cluster Quality

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

# Evaluation 1: Cluster Quality

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

# Evaluation 1: Cluster Quality

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

# Evaluation 1: Cluster Quality

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

# Evaluation 1: Cluster Quality

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

# Evaluation 1: Cluster Quality

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



# Evaluation 1: Cluster Quality

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

# Evaluation 1: Cluster Quality

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

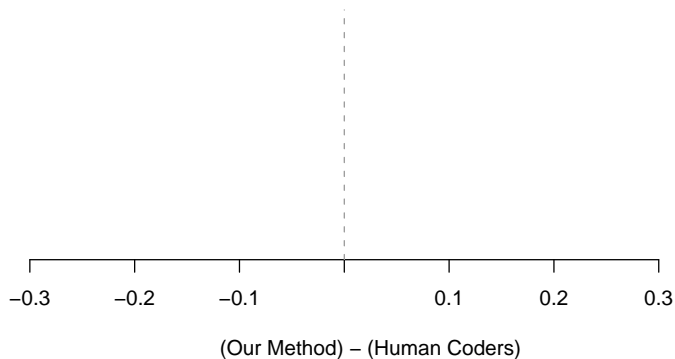
# Evaluation 1: Cluster Quality

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

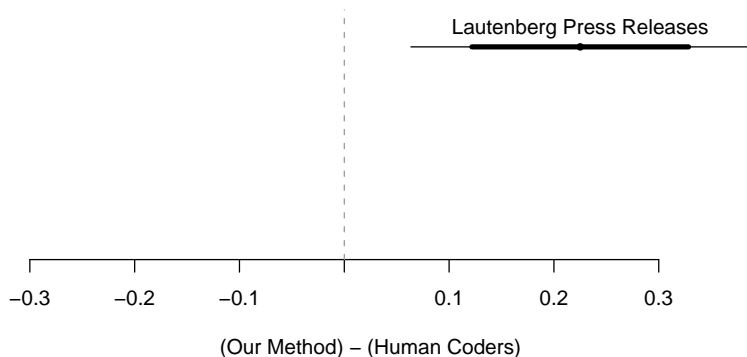
# Evaluation 1: Cluster Quality

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

# Evaluation 1: Cluster Quality

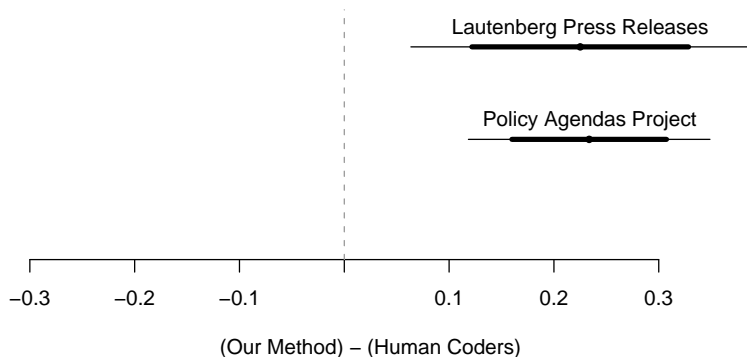


# Evaluation 1: Cluster Quality



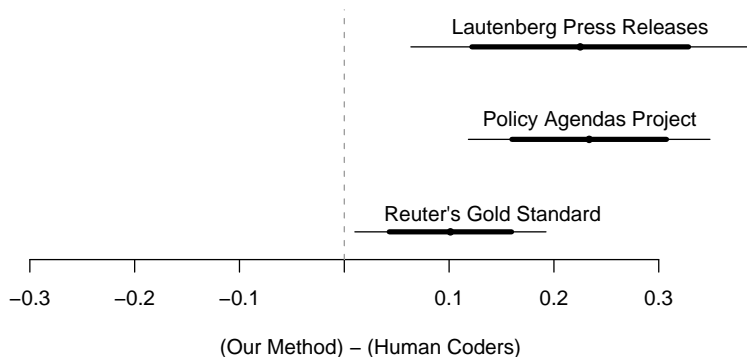
Lautenberg: 200 Senate Press Releases (appropriations, economy, education, tax, veterans, ...)

# Evaluation 1: Cluster Quality



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

# Evaluation 1: Cluster Quality



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



# Evaluation 2: More Informative Discoveries

## Evaluation 2: More Informative Discoveries

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

## Evaluation 2: More Informative Discoveries

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

## Evaluation 2: More Informative Discoveries

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

## Evaluation 2: More Informative Discoveries

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

## Evaluation 2: More Informative Discoveries

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

## Evaluation 2: More Informative Discoveries

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

## Evaluation 2: More Informative Discoveries

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



## Evaluation 2: More Informative Discoveries

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

## Evaluation 2: More Informative Discoveries

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

“Immigration”:

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

## Evaluation 2: More Informative Discoveries

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

“Immigration”:

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

“Genetic testing”:

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

# Evaluation 3: What Do Members of Congress Do?

# Evaluation 3: What Do Members of Congress Do?

- David Mayhew's (1974) famous typology

# Evaluation 3: What Do Members of Congress Do?

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

# Evaluation 3: What Do Members of Congress Do?

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

# Evaluation 3: What Do Members of Congress Do?

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



# Evaluation 3: What Do Members of Congress Do?

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

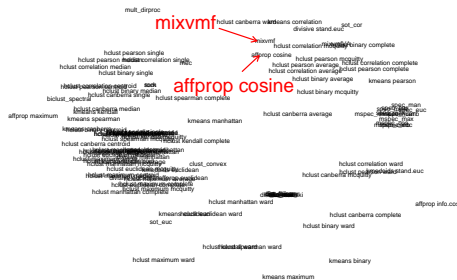
# Evaluation 3: What Do Members of Congress Do?

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





# Example Discovery



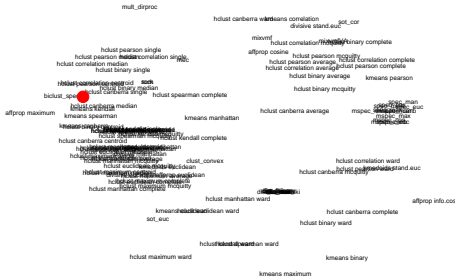
Red point: a **clustering** by Affinity Propagation-Cosine (Dueck and Frey 2007)

Close to:

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



# Example Discovery



Space between methods:

# Example Discovery



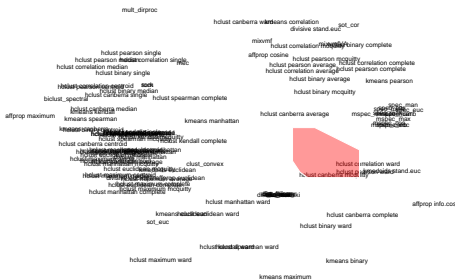
Space between methods:  
local cluster ensemble



# Example Discovery



# Example Discovery



Found a **region** with particularly insightful clusterings











# Example Discovery



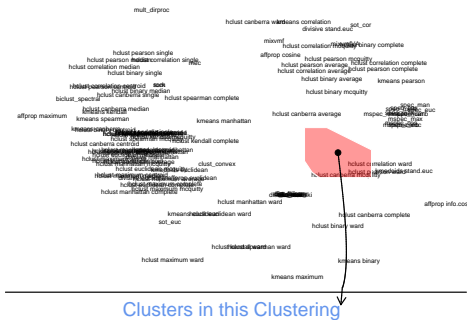
## Mixture:

- 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



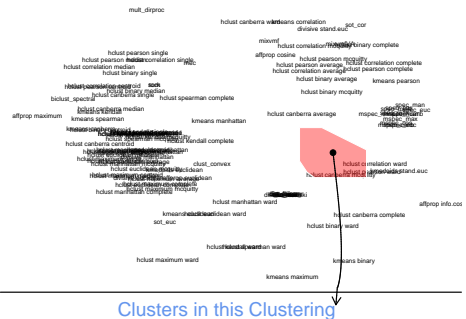


# Example Discovery





# Example Discovery



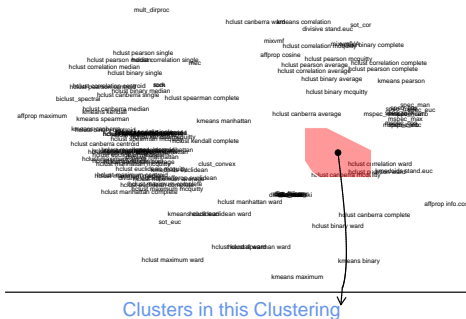
Credit Claiming  
Pork



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



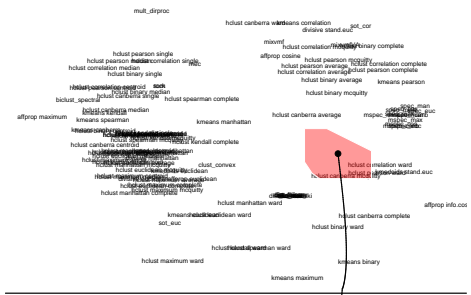
Credit Claiming  
Pork

Advertising

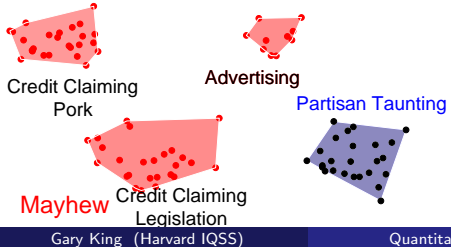
Mayhew  
Credit Claiming  
Legislation  
Gary King (Harvard IQSS)

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

# Example Discovery: Partisan Taunting

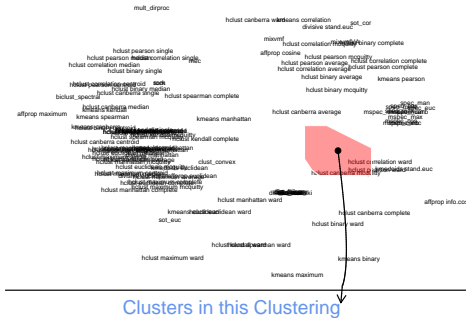


Clusters in this Clustering



Partisan Taunting:  
"Republicans Selling Out Nation  
on Chemical Plant Security"

# Example Discovery: Partisan Taunting



Credit Claiming  
Pork

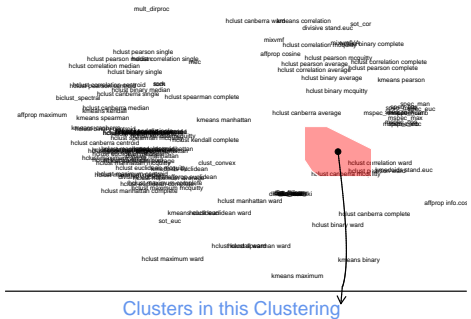
Advertising

Partisan Taunting

Mayhew  
Credit Claiming  
Legislation  
Gary King (Harvard IQSS)

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

# Example Discovery: Partisan Taunting



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



Credit Claiming  
Pork



Advertising

Partisan Taunting

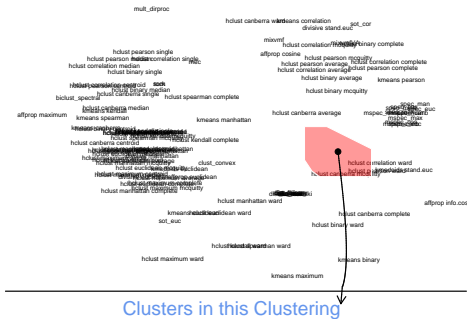


Mayhew  
Credit Claiming  
Legislation  
Gary King (Harvard IQSS)





# Example Discovery: Partisan Taunting



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

**Taunting ruins deliberation**



Credit Claiming  
Pork



Advertising

Partisan Taunting



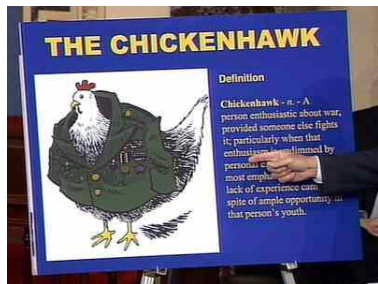
Mayhew  
Credit Claiming  
Legislation



Gary King (Harvard IQSS)

Quantitative Discovery

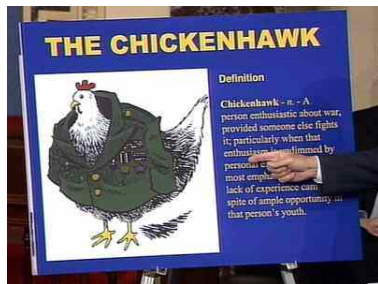
## Taunting ruins deliberation



Sen. Lautenberg  
on Senate Floor  
4/29/04

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

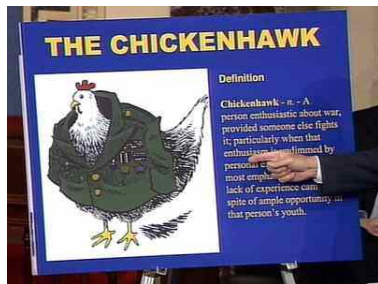
## Taunting ruins deliberation



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

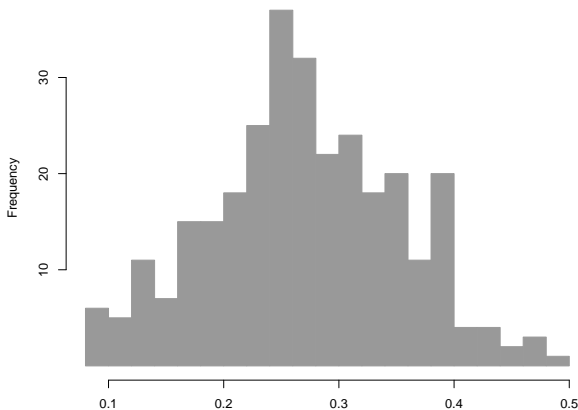
- Discovered using 200 press releases; 1 senator.
- Confirmed using 64,033 press releases; 301 senator-years.

# 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

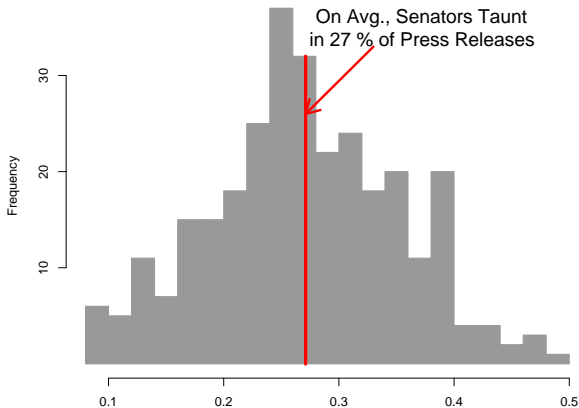
- Discovered using 200 press releases; 1 senator.
- Confirmed using 64,033 press releases; 301 senator-years.
- Apply supervised learning method: measure **proportion of press releases** a senator taunts other party



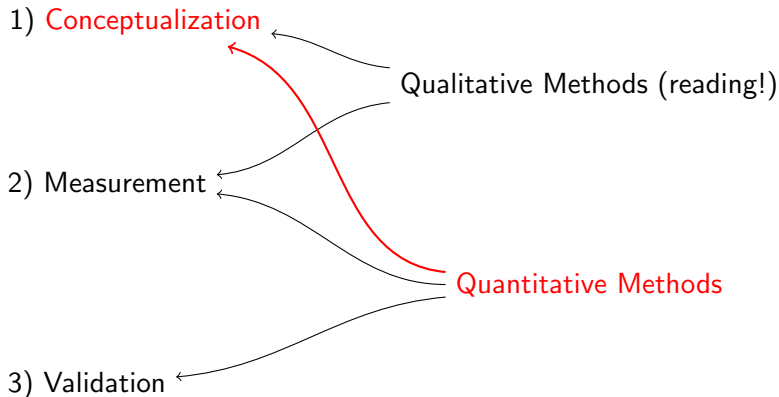


# Out of Sample Confirmation of Partisan Taunting

- 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

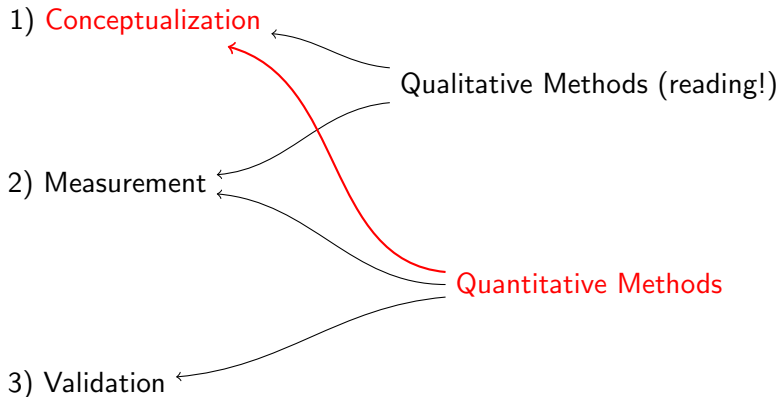


# Advancing the Objective of Discovery



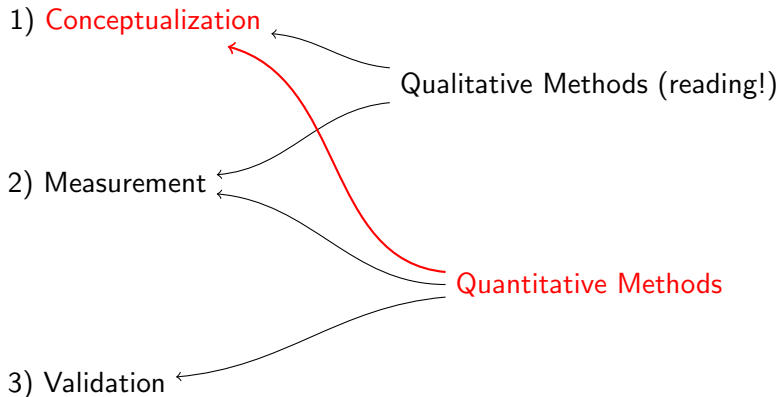
Quantitative methods for conceptualization: aiding **discovery**

# Advancing the Objective of Discovery



- Quantitative methods for conceptualization: aiding **discovery**
- Few formal methods designed explicitly for conceptualization

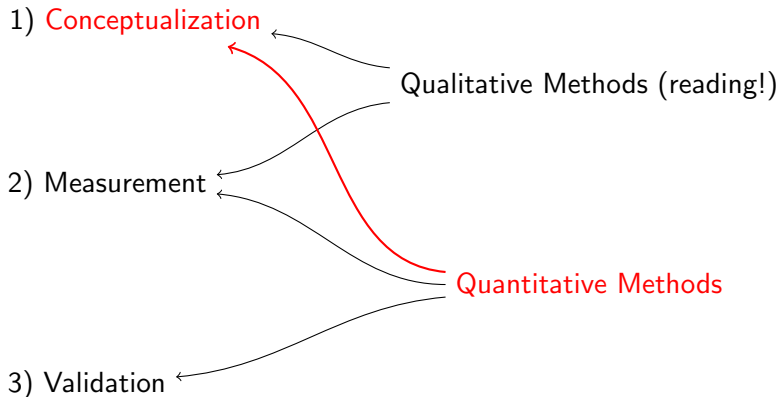
# Advancing the Objective of Discovery



Quantitative methods for conceptualization: aiding **discovery**

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

# Advancing the Objective of Discovery



Quantitative methods for conceptualization: aiding **discovery**

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

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

<http://GKing.Harvard.edu>