

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

Gary King

Institute for Quantitative Social Science  
Harvard University

Talk at BAE Systems, 9/9/2010

Joint work with Justin Grimmer (Harvard ↔ Stanford)

# 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*

# 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**
- No surprise: everyone's tried cluster analysis; very few are satisfied

# 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

# 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
- **The Question: How to organize all those clusterings?**

# Our Idea: Meaning Through Geography

Set of clusterings

# Our Idea: Meaning Through Geography

Set of clusterings  $\approx$

A list of unconnected addresses

wide at [SuperPages.com](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 Wob 02115..... 617 267-6483	
<b>Cartagena Lydia</b> 28 Sweet Box 02131..... 617 323-7639	<b>Faye &amp; Ricky</b> 207 Columbia Ave Bos 02136..... 617 437-7331	<b>Nicholas S F</b> 115 Randolph Ave Mil 02186..... 617 698-6307	
<b>Cartagena Avish</b> F Pleasant Rd 02139..... 617 442-9780	<b>Francis S</b> 134 Temple W Ave 02132..... 617 323-6781	<b>Nick 21 Farwell Box 02114..... 617 267-5222</b>	
<b>B Had 02134..... 617 361-5253</b>	<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 Howard Rd 02136..... 617 524-3078</b>	<b>Nicole..... 617 698-0713</b>	
<b>Luzmila</b> 124 Harvard Cam 02138..... 617 491-5621	<b>Fred W</b> 96 Newell Ave Mil 02186..... 617 698-1343	<b>Norman G</b> 38 Chickawhoh Dr 02125..... 617 822-1201	
<b>M 95 Howe Box 02132..... 617 323-9713</b>	<b>G &amp; B</b> 8 Vardon Dr 02134..... 617 436-8966	<b>P 40 Cranston Pl Bos 02115..... 617 457-4754</b>	
<b>Melvin</b> 503 Green Cam 02129..... 617 576-1061	<b>G T 27 Franklin Ave Som 02145..... 617 623-7121</b>	<b>P E 501 E South S Bos 02137..... 617 268-8213</b>	
<b>Carte Nicholas</b> 18 Appleton Boston 02114..... 617 695-6996	<b>Gayle</b> 25 Franklin St 02133..... 617 823-0322	<b>P E 14 Hutchings Box 02131..... 617 427-9170</b>	
<b>Cartier</b> 0 4 Bedford Box 02133..... 617 338-0219	<b>George</b> 125 Madison Bos 02134..... 617 367-9548	<b>P R 91 Boyer Ave 02138..... 617 968-8692</b>	
<b>Carten Thos J Sr &amp; Claire</b> 1 Franklin St Mil 02132..... 617 698-6163	<b>Carter Hillside Assoc</b> 107 S Street Bos 02111..... 617 456-1689	<b>Paul &amp; Constance</b> 114 Franklin St W Bos 02131..... 617 325-2036	
<b>17 445-5116</b>	<b>Carter Harry F</b> 30 Bayview Rd W Ave 02132..... 617 325-5465	<b>Paul E 501 E South S Bos 02137..... 617 268-4546</b>	
<b>17 822-2962</b>	<b>Carter Hide Co Inc</b> 26 Boston St 02111..... 617 542-7987	<b>Paul M 27 Crown St 02139..... 617 787-2115</b>	
<b>17 427-5712</b>	<b>A Nelson</b> 617 442-5230	<b>Carter Pile Driving Inc 27 Avenue G</b> Frankston 02702..... Wellesley Tpk 781.235-0488	
<b>17 569-2698</b>	<b>Carter Hilary 41 Harvey Cam 02148..... 617 876-2750</b>	<b>Carter Prudence</b> 34 Franklin Waterman 02127..... 617 393-3782	
<b>17 667-5190</b>	<b>Horace</b> 361 Walnut St Rosbury 02139..... 617 442-5307	<b>Prudence</b> 40 Franklin Waterman 02127..... 617 926-7063	
<b>17 569-1417</b>	<b>Howard Jr</b> 28 Nona Drive Box 02118..... 617 445-5532	<b>Roginald</b> 106 Brookview Dorchester 02122..... 617 541-2843	
<b>17 338-9110</b>	<b>J Dan..... 617 354-2658</b>	<b>Renee &amp; Andrew</b> 100 Walnut Bos 02138..... 617 720-3765	
<b>17 825-1953</b>	<b>J 31 Chatham Ave 02446..... 617 232-7990</b>	<b>Rice Dorel</b> 3400 Franklin Publishing 163 Main Wilmington 01887	
<b>17 296-1593</b>	<b>J 538 Harvard Bos 02446..... 617 730-9483</b>	<b>Ted Free-Dial '9' &amp; Thru..... 800 638-1671</b>	
<b>17 670-2078</b>	<b>J 775 The Pines West Rosbury 02132..... 617 323-5374</b>	<b>Carl Eric Industrial Prod 613 Main Wilmington</b> Ted Free-Dial '9' & Thru..... 800 619-7447	
<b>17 621-9001</b>	<b>J Breckin P Bn 02446..... 617 735-8787</b>	<b>Carl Free-Dial '9' &amp; Thru..... 800 648-7447</b>	
<b>17 296-4725</b>	<b>Carter J M</b> 1 Ipswich Pl Bos 02133..... 617 492-1214	<b>Carle</b> 113 Main Wilmington 0202	
<b>17 542-1521</b>	<b>B E 10 Graduate Ave Mil 02136..... 617 236-6329</b>	<b>Carl Free-Dial '9' &amp; Thru..... 978 988-7447</b>	
<b>17 364-5232</b>	<b>Carter Barbara L MD</b> Turfs New-England Medical Center Bos 02111	<b>Carl Richard A MD</b> 2079 Carverville Ave Brighton 02215..... 617 982-0836	
<b>17 541-5429</b>	<b>Carter Becky Jo 02134..... 617 523-4368</b>	<b>Carl Richard A MD</b> 130 Conant St Wob 02186..... 617 566-7293	
<b>17 739-2662</b>	<b>Bernard J</b> 122 Southside F Bus 02136..... 617 567-9430	<b>Carl Richard R J</b> 23 Mather St Som 02127..... 617 268-0448	
<b>17 879-0030</b>	<b>Bibbiah 25 Midway Dr 02134..... 617 298-8713</b>	<b>Carl Richard R J MD</b> 175 Rockdale Ave Cam 02142..... 617 864-1535	
<b>17 436-1513</b>	<b>Blair 28 Elmwood Pl 02138..... 617 367-9931</b>	<b>Roger 130 St Braughn Bos 02131..... 617 424-6148</b>	
<b>17 569-4119</b>	<b>Carl Broadcasting Co</b> 58 Park Pl Bos 02134..... 617 423-0210	<b>Roy 41 Concord Cam 02132..... 617 491-6115</b>	
<b>800 569-8782</b>	<b>Carl &amp; Susan Consultants Inc</b> 73 East St Cam 02541..... 617 225-0200	<b>Royce 18 Sanyday Cha 02129..... 617 241-0418</b>	
	<b>Carter C 200 Conantville Ave 02135..... 617 782-2118</b>		
	<b>C 218 Harvard Ave East Boston 02128..... 617 569-1545</b>		
	<b>C 109 Harvard Cam 02138..... 617 491-8522</b>		
	<b>C 101 Irving St Cambridge 02142..... 617 265-4932</b>		
	<b>C &amp; M 43 Bernham Jan 02136..... 617 524-9558</b>		



# 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 Ln Ipswich 01938	978 356-9960	
17 447-4101	Cartagena Lydia 28 Sweet Briar Rd 02131	617 323-7639	
90 257-9961	Cartagena Avish F Beach Rd 02139	617 442-9780	
17 566-1282	B Had 02136	617 361-5253	
17 364-5188	Lucille 124 Harvard Can 02139	617 491-5621	
361-0380	M 95 Howe Box 02136	617 323-9713	
17 566-4548	Melvin 503 Green Can 02139	617 576-1061	
17 628-8248	Carte Nicholas 18 Appleton Boston 02134	617 695-6996	
17 445-5116	Carters D & Claire 1 Furlow St Mt 02136	617 338-0219	
17 822-2962	Carte T & Kathleen 50 Thompson Ln Mt 02136	617 696-6919	
17 427-5712	A Weber A 200 Riverside Av Cambridge 02139	617 442-5230	
17 569-2698	A 21 Beulah Wy Haverhill 02119	617 442-1219	
17 667-5190	A M 250 Main St Av 02115	617 492-4174	
17 569-1417	Adams 301 Carter St Mt 02136	617 698-7074	
17 338-1107	Adams P 42 West St 02135	617 945-2711	
17 822-1959	Carte Anne MD 1161 Beacon St 02144	617 625-7623	
17 296-1193	Carte Anthony 971 Newbury Boston 02116	617 739-1022	
17 670-2078	B E 10 Gladstone Av Mt 02136	617 536-6229	
17 621-9001	Carte Barbara L MD Tufts New England Medical Center Box 02111	617 296-6911	
17 296-4725	Carte Becky MD 02134	617 636-0951	
17 542-1521	Bernard J 301 Ashdown E Mt 02136	617 523-4368	
17 364-5232	Bibb 25 Midway Dr 02136	617 567-9430	
17 541-5649	Billings 28 Newbury St 02138	617 298-8713	
17 739-2662	Carte Broadcasting Co 50 Park Pl Box 02136	617 367-9931	
17 879-0030	Carte C 2000 Cambridge St 02136	617 423-0210	
17 541-3948	Carte C 2000 Cambridge St 02136	617 225-0200	
17 436-1511	C 210 Townsend Av East Boston 02128	617 782-2118	
17 569-4119	C 109 Harvard Can 02136	617 569-1545	
909 569-8782	C & M 41 Northgate Jct 02134	617 491-8822	
	C & M 41 Northgate Jct 02134	617 524-9558	
	Carter F 514 Hicks Box 02131	617 327-1105	
	Faye & Ricky 20 Columbia Av Box 02136	617 437-7331	
	Francis S 134 Temple W Av 02132	617 323-6781	
	Franklin & Anne 705 Mt Auburn Can 02138	617 354-0798	
	Fred 41 Howard Av 02136	617 524-3078	
	Fred 76 Howley Av Mt 02136	617 698-1343	
	G & B 8 Vardon Box 02134	617 436-8906	
	G T 27 Fossil Av Mt 02145	617 623-7121	
	Gayle 25 Franklin St 02134	617 823-8322	
	Geo S 115 Main Mt Av 02136	617 522-3215	
	George 52 Madison Box 02134	617 367-9548	
	Carter Hillside Assoc 107 S Street Box 02111	617 456-1689	
	Carter Harry F 30 Bayview Rd W Av 02132	617 325-5465	
	Carter Hide Co Inc 140 Boston St W Av 02132	617 542-7987	
	Carter Hilary 41 Harvey Can 02148	617 876-2750	
	Horace 301 Walnut Av Haverhill 02119	617 442-5307	
	Howard Jr 28 New One Box 02118	617 445-5552	
	J Can 15 Chatham St 02144	617 324-2658	
	J 538 Harvard St 02146	617 232-7990	
	J 775 The Pine Way West Haverhill 02118	617 730-9483	
	Carter J Jacques MD 1 Breckinridge Pl Box 02144	617 735-8787	
	Carter J M 3410 Columbia Rd S Box 02137	617 464-1040	
	Carter J M Ornamental Ironworks 1000 Cambridge St 02139	617 436-5353	
	Carter J Neal Co 40 Newbury St 02138	617 442-1775	
	Carter James 1573 Cambridge St Can 02136	617 492-1214	
	James 102 Foster Av Haverhill 02118	617 739-2193	
	James 31 East Star Rd Cambridge 02141	617 876-8841	
	Jas L 34 Howley Rd Mt 02136	617 361-0773	
	Jane 14 Adams Rd Newton 02458	617 564-0435	
	John 1200 Cambridge St 02136	617 426-9094	
	John 11 Mansfield St 02134	617 987-2163	
	John 207 Summer St 02139	617 423-4334	
	John 40 Westmore St 02136	617 282-1235	
	Jane D 129 A Summit Av Box 02131	617 734-6109	
	J 290 Townsend Av East Boston 02128	617 265-8656	
	K 17 Concord Road 02123	617 282-1593	
	Carter Nellie E 323 Main St Av Box 02115	617 267-6483	
	Nicholas S F 115 Randolph Av Mt 02136	617 698-5307	
	Nick 21 Furlow Box 02116	617 267-5222	
	Nick & Debbi 136 Hermit Rd Newton 02459	617 527-0480	
	Norman G 38 Chickadee Dr 02126	617 822-1203	
	P 41 Eastwood Pl Box 02135	617 427-4754	
	P E 501 E South St Box 02137	617 268-8213	
	P L 44 Hutchings Box 02131	617 427-9170	
	P R 91 Boyer Can 02138	617 968-8692	
	Paul & Constance 114 Adams Av W Mt 02131	617 325-3034	
	Paul E 501 E South St Box 02137	617 268-4546	
	Paul M 27 Union St 02139	617 787-2115	
	Carter Pile Driving Inc 27 Beaver Ct Frankenm 02102	Wellesley Tpk-781.235-0488	
	Carter Prudence 40 Franklin Waterbury 02172	617 393-3782	
	Prudence 40 Franklin Waterbury 02172	617 926-7063	
	Reginald 100 Broadmarch Chester 02124	617 541-2843	
	Renée & Andrew 30 Walnut St 02138	617 720-3765	
	Carter Rice David Building Division 163 Main Wilmington 01887 Toll Free 800 714 7136	800 638-1671	
	Robt 4100 Cambridge St 02139	800 619-7447	
	Toll Free 800 714 7136	800 619-7447	
	Toll Free 800 714 7136	800 648-7447	
	Wilmington 413 Main Wilmington 01887	978 988-7447	
	Ingalls Centre 163 Main Wilmington 01887	800 638-1673	
	Carter Richard 2079 Cambridge Av Brighton 02116	617 987-0836	
	Richard A 97 W Vernon St 02136	617 566-7293	
	Carter Richard A 120 Cambridge St 02136	617 267-0710	
	Carter Richard K 123 Merwin St Box 02137	617 268-0468	
	Robert L 175 Newbury Av Can 02136	617 864-1535	
	Rose 120 A Summit Av Box 02131	617 424-6148	
	Royce & Andrew 180 Broadway Av 02129	617 491-6115	
	Royce 18 Broadway Av 02129	617 241-9418	



$\approx$  We develop a (conceptual) geography of clusterings



# A New Strategy

Make it easy to choose best clustering from millions of choices

# A New Strategy

Make it easy to choose best clustering from millions of choices

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

# A New Strategy

Make it easy to choose best clustering from millions of choices

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

# A New Strategy

Make it easy to choose best clustering from millions of choices

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

# A New Strategy

Make it easy to choose best clustering from millions of choices

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

# A New Strategy

Make it easy to choose best clustering from millions of choices

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

# A New Strategy

Make it easy to choose best clustering from millions of choices

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

# A New Strategy

Make it easy to choose best clustering from millions of choices

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





# 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

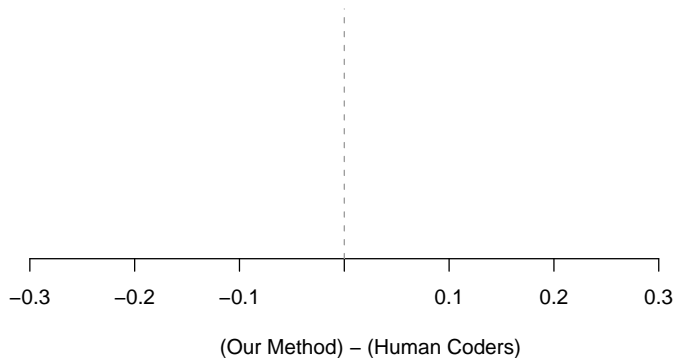
- **Experimental Design to Assess Cluster Quality**

- 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**

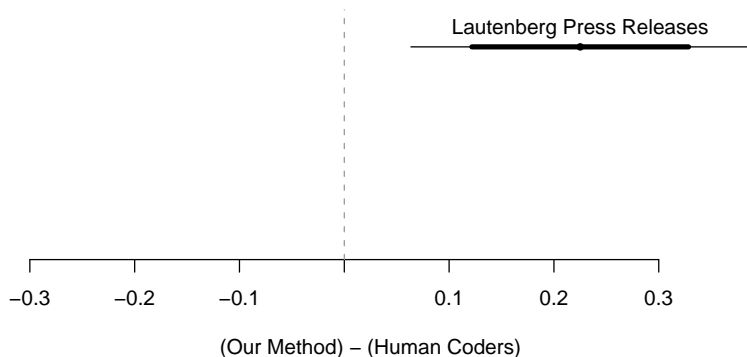
# Evaluation 1: Cluster Quality

- **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

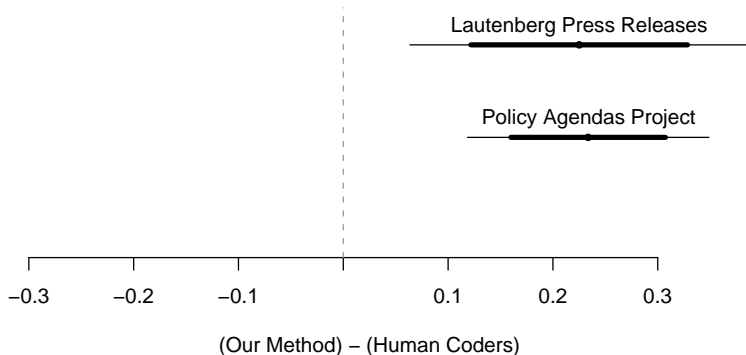


# Evaluation 1: Cluster Quality



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

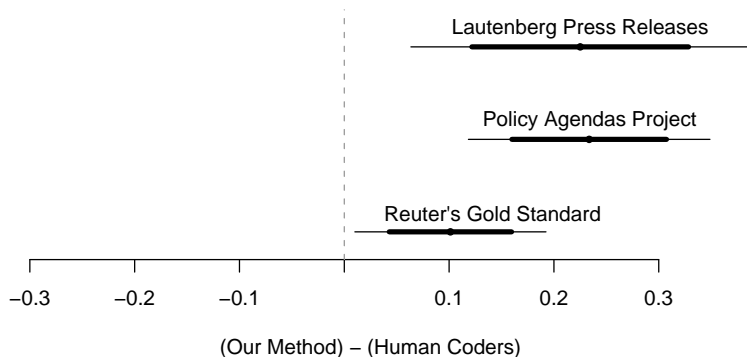
# Evaluation 1: Cluster Quality



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





























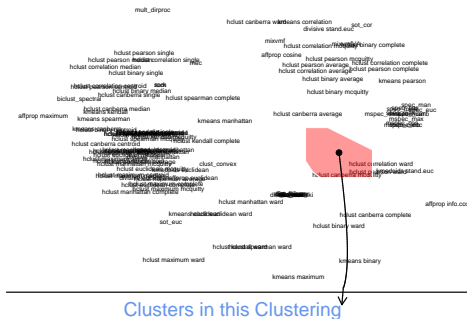








# Example Discovery

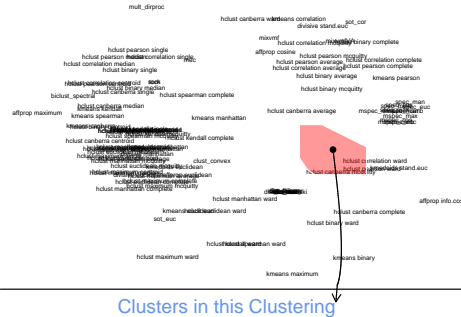


Credit Claiming  
Pork

## Credit Claiming, Pork:

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

# Example Discovery



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



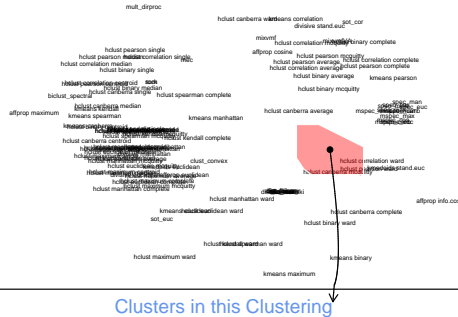
Credit Claiming  
Pork



Mayhew Credit Claiming  
Legislation

Gary King (Harvard IQSS)

# Example Discovery



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

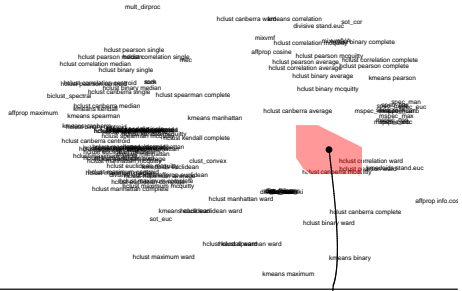
Credit Claiming  
 Pork

Advertising

Mayhew  
 Credit Claiming  
 Legislation

Gary King (Harvard IQSS)

# Example Discovery: Partisan Taunting



Clusters in this Clustering

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



Credit Claiming  
 Pork



Advertising

Partisan Taunting



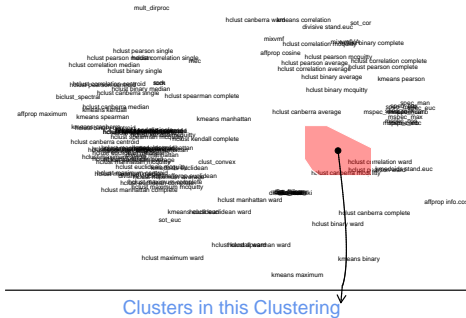
Mayhew Credit Claiming  
 Legislation



Gary King (Harvard IQSS)

Quantitative Discovery

# Example Discovery: Partisan Taunting



Credit Claiming  
Pork

Advertising

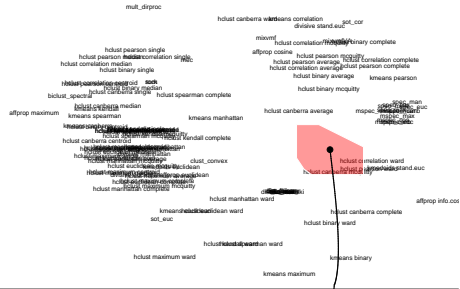
Partisan Taunting

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



Clusters in this Clustering



Credit Claiming  
Pork



Advertising



Mayhew  
Credit Claiming  
Legislation

Gary King (Harvard IQSS)

Partisan Taunting

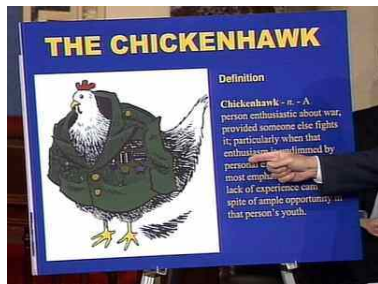


Quantitative Discovery

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



## Taunting ruins deliberation

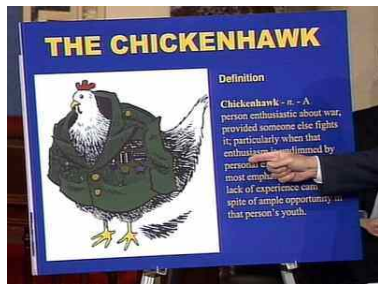


Sen. Lautenberg  
on Senate Floor  
4/29/04

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

# In Sample Illustration of Partisan Taunting

## Taunting ruins deliberation

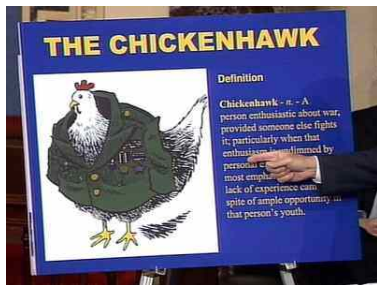


Sen. Lautenberg  
on Senate Floor  
4/29/04

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

# In Sample Illustration of Partisan Taunting

## Taunting ruins deliberation



Sen. Lautenberg  
on Senate Floor  
4/29/04

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

# Out of Sample Confirmation of Partisan Taunting

- Discovered using 200 press releases; 1 senator.

# Out of Sample Confirmation of Partisan Taunting

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

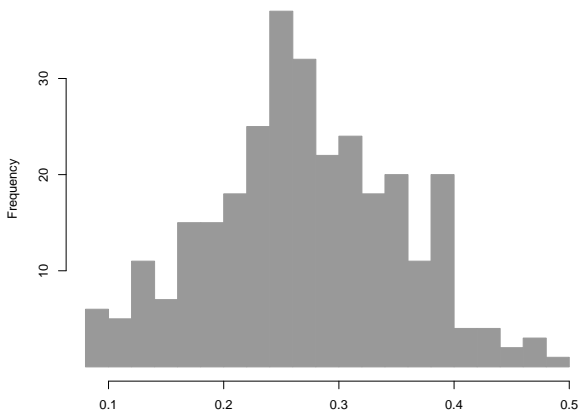
# Out of Sample Confirmation of Partisan Taunting

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



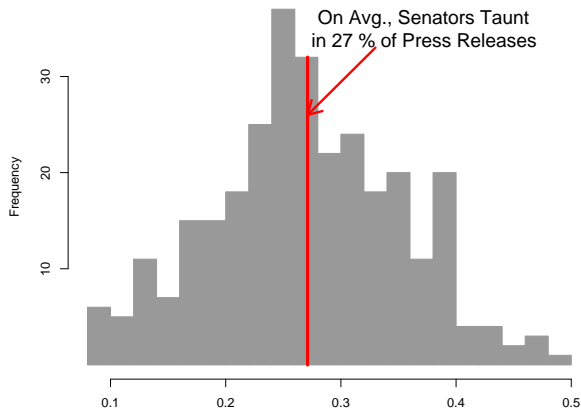
# Out of Sample Confirmation of Partisan Taunting

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

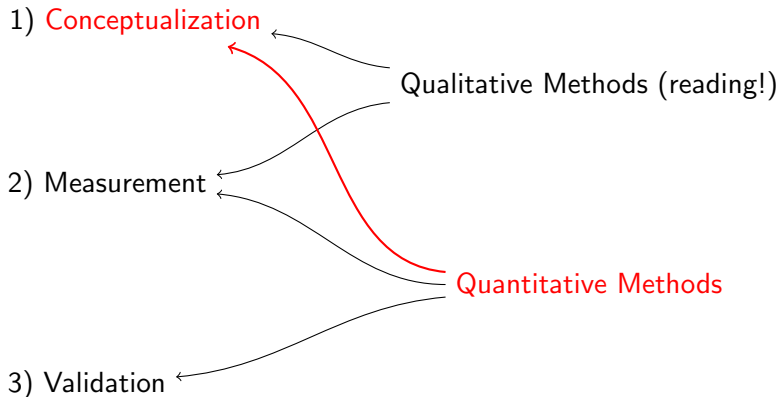


# Out of Sample Confirmation of Partisan Taunting

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

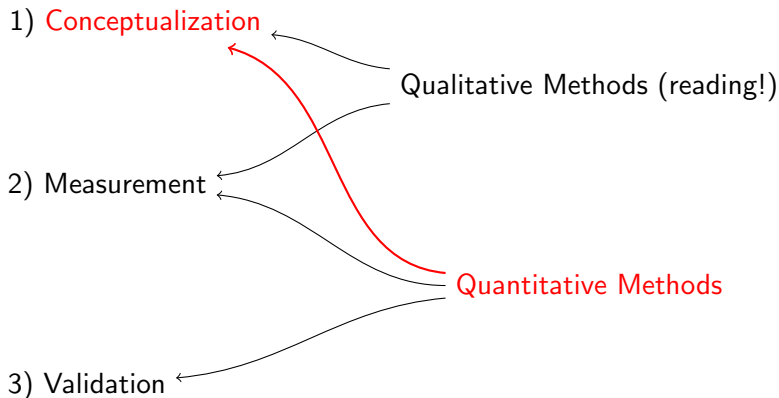


# Advancing the Objective of Discovery



Quantitative methods for conceptualization: aiding **discovery**

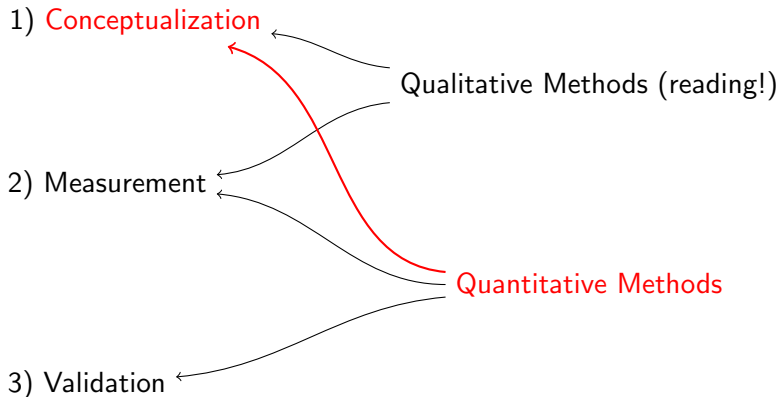
# Advancing the Objective of Discovery



Quantitative methods for conceptualization: aiding **discovery**

- Few formal methods designed explicitly for conceptualization

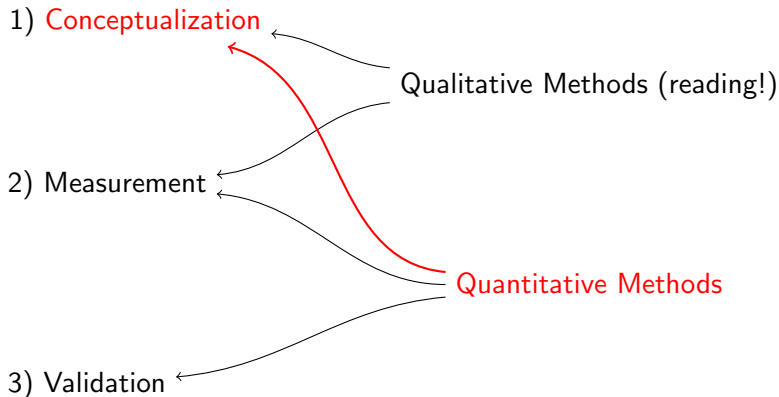
# Advancing the Objective of Discovery



Quantitative methods for conceptualization: aiding **discovery**

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

# Advancing the Objective of Discovery



Quantitative methods for conceptualization: aiding **discovery**

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

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

<http://GKing.Harvard.edu>