

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

Gary King
Institute for Quantitative Social Science
Harvard University

joint work with
Justin Grimmer (Harvard University)

(talk at the Northeast Methodology Program, New York University, 4/17/09)

The Problem: Discovery from Unstructured Text

- Examples: scholarly literature, news stories, medical information, blog posts, comments, product reviews, emails, social media updates, audio-to-text summaries, speeches, press releases, legal decisions, etc.
- 10 minutes of worldwide email = 1 LOC equivalent
- An essential part of discovery is **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**: discovery through (1) classification and (2) simultaneously inventing a classification scheme
- (We analyze text; our methods apply more generally)

Why Johnny Can't Classify (Optimally)

- $Bell(n)$ = 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!
- That we think of all this as astonishing ... is astonishing

Why HAL Can't Classify Either

- The Goal — an optimal application-independent cluster analysis method — is mathematically impossible:
 - **Ugly duckling theorem**: every pair of documents are equally similar
↔ every partition of documents is equally similar
 - **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, **who knows?!**
 - The literature: **little guidance on when methods apply**
 - **Deriving such guidance: difficult or impossible**
 - (Perhaps true by definition in unsupervised learning: If we knew the DGP, we wouldn't be at the discovery stage.)

If Ex Ante doesn't work, try Ex Post

- **Methods and substance must be connected** (no free lunch theorem)
- The usual approach fails: can't do it by understanding the model
- We do it **ex post** (by cherry-picking results)
 - For discovery (our goal): No problem
 - For estimation & confirmation: more difficult or biased
- Complicated concepts are easier to define ex post:
 - “I know it when I see it” (Justice Stewart's definition of obscenity)
 - Anchoring Vignettes (on defining concepts by example)
- **But how to choose from an enormous list of clusterings?**

Our Idea: Meaning Through Geography

wide at SuperPages.com

195 Car

17 566-1282	Cartage New England Inc 36 Abbott St Ipswich 01938... 978 356-9960	Carter F 34 Inwood Ave 02113... 617 327-1105	Carter Nella E 321 Main St 02113... 617 267-6483
17 447-4101	Cartagema Lydia 30 Liberty St 02113... 617 323-7639	Faye & Ricky 707 Chatham Ave 02116... 617 437-7331	Nicholas S F 125 Northrop Ave 02186... 617 498-5307
800 257-9981	Cartagema Aethi 9 Beach St 02119... 617 442-9780	Franklin & Anne 271 Mt Auburn Ave 02138... 617 354-4708	Nick & Debbie 134 Harvard Rd Newton 02459... 617 527-0480
17 566-1282	Bart 02138 Jessica 30 Decatur Cir 02129... 617 241-0152	Fred 47 Woodlawn Ave 02138... 617 434-3078	Nicole 134 Harvard Rd Newton 02459... 617 498-0713
17 364-5188	Lucilla 174 Harvard Cir 02139... 617 491-5621	Fred W Hinchey III MD 02138... 617 698-1343	Norman G 38 Chickatabuck Dr 02132... 617 822-1203
361-0380	M 92 Tower St 02138... 617 303-9713	G T 27 Franklin Ave 02146... 617 623-7121	P 34 Croswell Pl 02133... 617 427-4754
17 566-4548	Merlin 321 Green Can 02134... 617 576-1061	Gayle 25 Front St 02134... 617 825-0322	P E 161 E South St 02127... 617 248-0213
Carle Nicholas	Carte 11000 Carton Thos Jr & Co 02136... 617 695-6996	Geo S 155 Mount Mt St 02133... 617 522-3215	P L 44 Hollings Ave 02114... 617 427-9170
17 628-8248	Cartena O 4 Miller St 02118... 617 338-8219	George 125 Hazden Ave 02114... 617 367-9548	P R Jr 250 Main St 02114... 617 983-8692
17 445-5116	Carten Thos Jr & Co 02136... 617 698-6163	Carter Halliday Associate 380 S Tower St 02113... 617 456-1689	Paul & Constantine 114 Ansonia Av W 02132... 617 325-2036
17 822-2982	Thomas & Kathleen 50 Thompson Ln 02136... 617 696-6919	Carter Harry E 100 Irving St 02113... 617 325-5465	Paul E 303 E South St 02127... 617 268-0546
17 427-5712	A Industry Carter A 02111... 617 327-2257	Carter Hide Co Inc 146 Summer St 02133... 617 542-7987	Paul M 37 Union St 02135... 617 787-2115
17 569-2608	A 31 Redburn Way 02139... 617 442-1219	Carter Hilary 61 Harvey Can 02186... 617 076-2750	Prudence 40 Francis Walden Dr 02127... 617 393-3782
17 667-5190	A M 255 Montross Ave 02135... 617 266-7153	Horace 362 Walnut Av 02139... 617 442-5307	Prudence 40 Francis Walden Dr 02127... 617 986-7063
17 569-1417	Adams 361 Centre St 02138... 617 698-9074	Howard Jr 30 Main Dr 02139... 617 445-5552	Reynold 336 Wrentham Dr 02132... 617 541-2843
17 338-8110	Alice 108 Elmwood St 02135... 617 425-0193	J Canard 617 354-2688	Renée & Andrew 30 Walnut St 02138... 617 720-3765
17 825-9195	Alice 45 Market Cambridge 02139... 617 945-2711	J 25 Chatham St 02144... 617 232-7990	Renee & Andrew 30 Walnut St 02138... 617 720-3765
17 296-1593	Andrew F 22 Wood Av 02143... 617 625-7623	J 518 Harvard St 02146... 617 730-9483	Renee & Andrew 30 Walnut St 02138... 617 720-3765
17 670-2078	Andrew IV 02144... 617 739-1022	J 775 The Pines Wood 02121... 617 323-5574	Renee & Andrew 30 Walnut St 02138... 617 720-3765
17 623-9001	Carter Arthens 271 Beulah Boston 02114... 617 536-6329	Carter J Jacques MD 1101 Beacon St 02146... 617 739-1022	Renee & Andrew 30 Walnut St 02138... 617 720-3765
17 296-4723	B E 48 Gloucester Av 02126... 617 296-6911	Carter J M 1433 Columbia St S 02137... 617 464-1040	Renee & Andrew 30 Walnut St 02138... 617 720-3765
17 296-4723	Carter Barbara L MD 100 New England Medical Center 02113	Carter J M Ornamental Ironworks 1433 Columbia St S 02137... 617 464-1040	Renee & Andrew 30 Walnut St 02138... 617 720-3765
17 542-1521	Cel... 617 636-0051	Carter J M Ornamental Ironworks 1433 Columbia St S 02137... 617 464-1040	Renee & Andrew 30 Walnut St 02138... 617 720-3765
17 364-5232	Carter Becky Jo 02134... 617 523-4568	Carter J M Ornamental Ironworks 1433 Columbia St S 02137... 617 464-1040	Renee & Andrew 30 Walnut St 02138... 617 720-3765
17 541-5649	Carter Bernard J 121 Dan St 02138... 617 567-3430	Carter J M Ornamental Ironworks 1433 Columbia St S 02137... 617 464-1040	Renee & Andrew 30 Walnut St 02138... 617 720-3765
17 739-2662	Bobbi 20 Wilmot Dr 02124... 617 290-8713	Carter J M Ornamental Ironworks 1433 Columbia St S 02137... 617 464-1040	Renee & Andrew 30 Walnut St 02138... 617 720-3765
17 879-0030	Blake 20 Wilmot Dr 02124... 617 290-8713	Carter J M Ornamental Ironworks 1433 Columbia St S 02137... 617 464-1040	Renee & Andrew 30 Walnut St 02138... 617 720-3765
17 541-3948	Blake 20 Wilmot Dr 02124... 617 290-8713	Carter J M Ornamental Ironworks 1433 Columbia St S 02137... 617 464-1040	Renee & Andrew 30 Walnut St 02138... 617 720-3765
17 436-1513	Blake 20 Wilmot Dr 02124... 617 290-8713	Carter J M Ornamental Ironworks 1433 Columbia St S 02137... 617 464-1040	Renee & Andrew 30 Walnut St 02138... 617 720-3765
17 569-4119	Blake 20 Wilmot Dr 02124... 617 290-8713	Carter J M Ornamental Ironworks 1433 Columbia St S 02137... 617 464-1040	Renee & Andrew 30 Walnut St 02138... 617 720-3765
800 569-8782	Blake 20 Wilmot Dr 02124... 617 290-8713	Carter J M Ornamental Ironworks 1433 Columbia St S 02137... 617 464-1040	Renee & Andrew 30 Walnut St 02138... 617 720-3765



↪ We develop a geography of clusterings

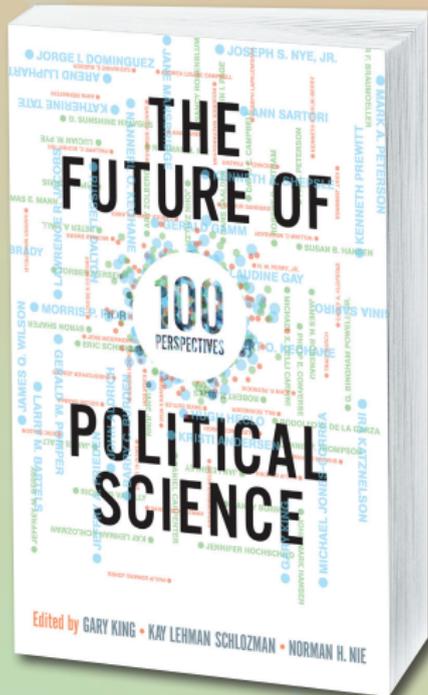
A New Strategy

- 1 Code text as numbers (in one or more of several ways)
- 2 Apply all existing clustering methods (that have been used by at least one person other than the author) to the data — each representing different substantive assumptions (<15 mins)
- 3 Develop an application-independent distance metric between clusterings
- 4 Create a metric space of clusterings, and a 2D projection
- 5 Introduce the local cluster ensemble to summarize any point, including points with no existing clustering
- 6 Propose a new animated visualization: use the local cluster ensemble to explore the space of clusterings (smoothly morphing from one into others)

↪ meaning revealed through a geography of clusterings

Application-Independent Distance Metric: Axioms

- 1 Clusterings with more **pairwise document agreements** are closer (we prove: pairwise agreements encompass 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 $\log(\text{num clusters})$
- ↪ **Only one measure satisfies all three** (the “variation of information”)



Available March 2009: 304pp
Pb: 978-0-415-99701-0: **\$24.95**
www.routledge.com/politics

THE FUTURE OF POLITICAL SCIENCE

100 Perspectives

Edited by Gary King, Harvard University, Kay Lehman Schlozman, Boston College
and Norman H. Nie, Stanford University

“The list of authors in *The Future of Political Science* is a ‘who’s who’ of political science. As I was reading it, I came to think of it as a platter of tasty hors d’oeuvres. It hooked me thoroughly.”

—Peter Kingstone, University of Connecticut

“In this one-of-a-kind collection, an eclectic set of contributors offer short but forceful forecasts about the future of the discipline. The resulting assortment is captivating, consistently thought-provoking, often intriguing, and sure to spur discussion and debate.”

—Wendy K. Tam Cho, University of Illinois at Urbana-Champaign

“King, Schlozman, and Nie have created a visionary and stimulating volume. The organization of the essays strikes me as nothing less than brilliant. . . It is truly a joy to read.”

—Lawrence C. Dodd, Manning J. Dauer Eminent Scholar in Political Science,
University of Florida

 **Routledge**
Taylor & Francis Group
an **informa** business

Evaluators' Rate Machine Choices Better Than Their Own

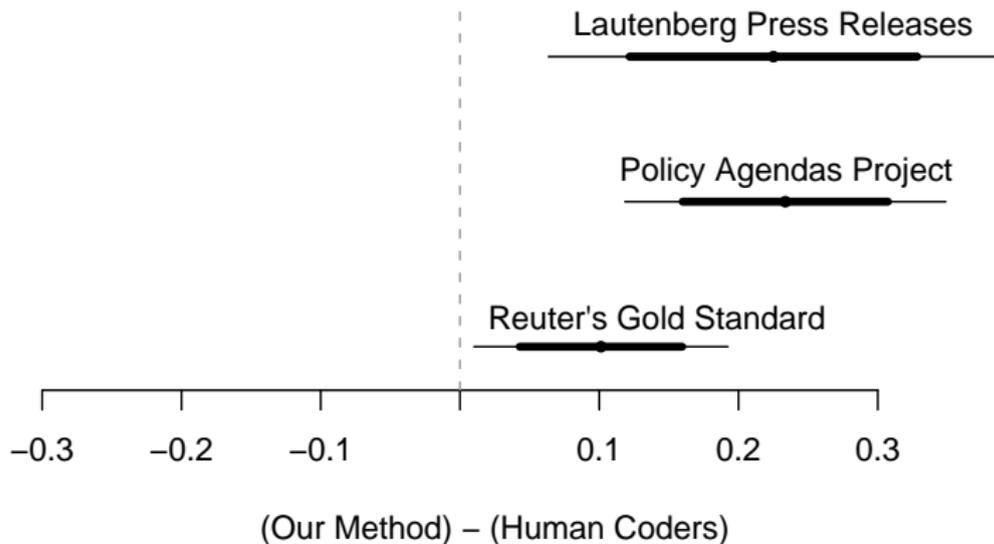
- Scale: (1) unrelated, (2) loosely related, or (3) closely related
- Table reports: mean(scale)

Pairs from	Overall Mean	Evaluator 1	Evaluator 2
Random Selection	1.38	1.16	1.60
Hand-Coded Clusters	1.58	1.48	1.68
Hand-Coding	2.06	1.88	2.24
Machine	2.24	2.08	2.40

p.s. The hand-coders did the evaluation!

Cluster Quality Experiments

Scale: mean(within clusters) – mean(between clusters)



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

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

Reuter's: financial news (trade, earnings, copper, gold, coffee, ...): "gold

What do Members of Congress Do?

Substantive example of a finding, using our approach

- David Mayhew's (1974) famous typology
 - ① Advertising
 - ② Credit Claiming
 - ③ Position Taking
- We find one more: **Partisan Taunting**
 - "Senator Lautenberg Blasts Republicans as 'Chicken Hawks' " [Government Oversight]
 - "The Intolerance and discrimination from the Bush administration against gay and lesbian Americans is astounding" [Civil Rights]
 - "John Kerry had enough conviction to sign up for the military during wartime, unlike the Vice President, who had a deep conviction to avoid military service" [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]
 - ↪ **Perhaps this is what it means to be a member of a political party in the U.S.?**

More Informative Discoveries

- Found 2 scholars analyzing lots of textual data for their work
- For each: created 2 clusterings from each of 3 methods, including ours
- Created info packet on each clustering (for each cluster: exemplar document, automated content summary)
- Asked for $\binom{6}{2}=15$ pairwise comparisons
- Both cases a Condorcet winner:

“Immigration”:

Our Method 1 → vMF 1 → vMF 2 → Our Method 2 → K-Means 1 → K-Means 2

“Genetic testing”:

Our Method 1 → {Our Method 2, K-Means 1, K-means 2} → Dir Proc. 1 → Dir Proc. 2

- **Intended contributions:**

- An encompassing cluster analytic approach for discovery
- A new approach to evaluating results in unsupervised learning
- Especially useful for the ongoing spectacular increase in the production and availability of unstructured text

- **Future research:**

- Advancing our approach: (1) $>2D$ exploration, (2) alternative visualizations of the space of clusterings, (3) including more methods
- Evaluating new individual methods: (1) distance from existing methods and their averages, (2) usefulness of discoveries in given data sets.

For more information:

<http://GKing.Harvard.edu>