

# Discovery

Gary King  
Institute for Quantitative Social Science  
Harvard University

covering joint work with  
Justin Grimmer (Harvard) and Eleanor Powell (Harvard ↔ Yale)

(talk @ Institute for Qualitative and Multimethod Research, Syracuse University, 5/26/09)

- Gary King and Eleanor Neff Powell. 2008. “How Not to Lie Without Statistics”
- Justin Grimmer and Gary King. 2009. “Quantitative Discovery of Qualitative Information: A General Purpose Document Clustering Methodology”

<http://GKing.harvard.edu>

# Is Qualitative or Quantitative Research Better?

This question has an answer (hard to know from the qualitative methods literature!)

## If insufficient information is quantified, qualitative judgment wins

- The vast majority of human decisions are qualitative
- If not, we wouldn't have survived as a species

## If sufficient information is quantified, statistics wins

- The information required: shockingly little
- 6 crude variables beat 83 law professors (Martin et al., 2004)
- Simple quantitative models forecast elections better than expert pundits (Gelman and King, 1993; Campbell, 2005)
- 284 experts forecast the political future with “less skill than simple extrapolation” (Tetlock, 2005)
- statistical model out-performs physicians in determining cause-specific mortality rates (King and Lu, 2008)
- Hundreds of head-to-head contests, mostly with the same conclusion (Meehl, 1954; Grove, 2005)
- The **march of quantification** across fields of academia, professions, commerce, sports, etc. (Moneyball, SuperCrunchers, Numerati)

## 3 Steps

- 1 **Conceptualization** (e.g., a categorization scheme)
- 2 **Measurement** (e.g., classifying objects into categories)
- 3 **Verification** (e.g., testing a hypothesis that could be wrong)

## Styles & Approaches (the cause of many misunderstandings)

- **Quantitative**: focus on measurement and verification  $\rightsquigarrow$  assumes & underplays conceptualization
- **Qualitative**: focus on (& iterate between) conceptualization and measurement  $\rightsquigarrow$  ignores or underplays verification

## Proposed Solutions

- Do both qualitative & quantitative analysis, separately
- Atheoretical unsupervised learning algorithms
- $\rightsquigarrow$  Integrated quant/qual: **Computer-assisted qualitative analysis**

# 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
- **Qualitative-only approaches are hopeless**
- 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:
  - **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 qualitative choice)
  - 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?**





# A New Strategy

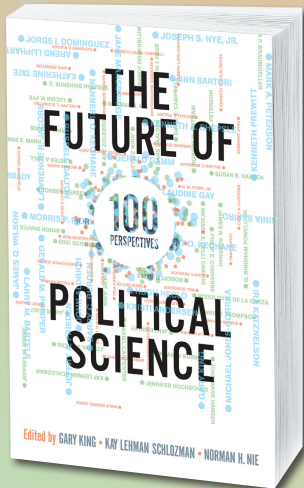
- 1 **Code text as numbers** (in one *or more* of several ways)
- 2 **Apply all existing clustering methods** (that have been used by at least one person other than the author) to the data — each representing different substantive assumptions (<15 mins)
- 3 Develop an **application-independent distance** 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](http://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

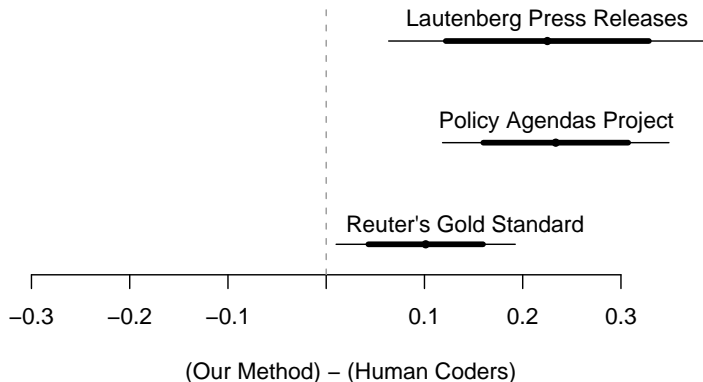
- Scale: (1) unrelated, (2) loosely related, or (3) closely related
- Table reports: mean(scale)

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
<b>Machine</b>	<b>2.24</b>	<b>2.08</b>	<b>2.40</b>

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

# Cluster Quality Experiments

Scale:  $\text{mean}(\text{within clusters}) - \text{mean}(\text{between clusters})$



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 scopes trial took place in 1925. Sadly, President Bush’s veto today shows that we haven’t progressed much since then.” [Healthcare]
  - “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]
  - ↪ **Is this what it means to be a member of a political party?**



# 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
- **Multiple approaches: Integrated, not separate**

For more information:

<http://GKing.Harvard.edu>