# Computer-Assisted Clustering and Conceptualization from Unstructured Text

Gary King

Institute for Quantitative Social Science
Harvard University

talk at University of Chicago, Computation Institute, 5/9/2011

---

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

# A Method for Computer Assisted Conceptualization

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

- Cluster Analysis: simultaneously (1) invents categories and (2) assigns documents to categories

# A Method for Computer Assisted Conceptualization

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

- Cluster Analysis: simultaneously (1) invents categories and (2) assigns documents to categories

- We focus on unstructured text; methods apply more broadly.

# A Method for Computer Assisted Conceptualization

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

- Cluster Analysis: simultaneously (1) invents categories and (2) assigns documents to categories

- We focus on unstructured text; methods apply more broadly.

- Main goal: Switch from Fully Automated to Computer Assisted

# What's Hard about Clustering?

# What's Hard about Clustering?
(aka Why Johnny Can't Classify)

- Clustering seems easy; its not!

# What's Hard about Clustering?
(aka Why Johnny Can't Classify)

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

# What's Hard about Clustering?
## (aka Why Johnny Can't Classify)

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

# What's Hard about Clustering?
## (aka Why Johnny Can't Classify)

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

# What's Hard about Clustering?
(aka Why Johnny Can't Classify)

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

# What's Hard about Clustering?
(aka Why Johnny Can't Classify)

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

# What's Hard about Clustering?
(aka Why Johnny Can't Classify)

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

# What's Hard about Clustering?
(aka Why Johnny Can't Classify)

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

# What's Hard about Clustering?
(aka Why Johnny Can't Classify)

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

# The Problem with Fully Automated Clustering

- The (Impossible) Goal: optimal, fully automated, application-independent cluster analysis

# The Problem with Fully Automated Clustering

- The (Impossible) Goal: optimal, fully automated, application-independent cluster analysis
- No free lunch theorem: every possible clustering method performs equally well on average over all possible substantive applications

# The Problem with Fully Automated Clustering

- The (Impossible) Goal: optimal, fully automated, application-independent cluster analysis
- No free lunch theorem: every possible clustering method performs equally well on average over all possible substantive applications
- Existing methods:

# The Problem with Fully Automated Clustering

- The (Impossible) Goal: optimal, fully automated, application-independent cluster analysis
- No free lunch theorem: every possible clustering method performs equally well on average over all possible substantive applications
- Existing methods:
  - Many choices: model-based, subspace, spectral, grid-based, graph-based, fuzzy $k$-modes, affinity propagation, self-organizing maps,...

# The Problem with Fully Automated Clustering

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

# The Problem with Fully Automated Clustering

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

# The Problem with Fully Automated Clustering

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

# The Problem with Fully Automated Clustering

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

# The Problem with Fully Automated Clustering

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

# The Problem with Fully Automated Clustering

- **The (Impossible) Goal:** optimal, fully automated, application-independent cluster analysis
- **No free lunch theorem:** every possible clustering method performs equally well on average over all possible substantive applications
- Existing methods:
  - **Many choices:** model-based, subspace, spectral, grid-based, graph-based, fuzzy $k$-modes, affinity propagation, self-organizing maps,...
  - **Well-defined** statistical, data analytic, or machine learning foundations
  - How to add substantive knowledge: With few exceptions, unclear
  - The literature: little guidance on when methods apply
  - Deriving such guidance: difficult or impossible
- **Deep problem: full automation requires more information**
- No surprise: everyone's tried cluster analysis; very few are satisfied

# Switch from Fully Automated to Computer Assisted

- Fully Automated Clustering may succeed sometimes, but fails in general: too hard to understand when each model applies

# Switch from Fully Automated to Computer Assisted

- Fully Automated Clustering may succeed sometimes, but fails in general: too hard to understand when each model applies
- An alternative: Computer-Assisted Clustering

# Switch from Fully Automated to Computer Assisted

- Fully Automated Clustering may succeed sometimes, but fails in general: too hard to understand when each model applies
- An alternative: Computer-Assisted Clustering
  - Easy in theory: list all clusterings; choose the best

# Switch from Fully Automated to Computer Assisted

- Fully Automated Clustering may succeed sometimes, but fails in general: too hard to understand when each model applies
- An alternative: Computer-Assisted Clustering
  - Easy in theory: list all clusterings; choose the best
  - Impossible in practice: Too hard for us mere humans!

# Switch from Fully Automated to Computer Assisted

- Fully Automated Clustering may succeed sometimes, but fails in general: too hard to understand when each model applies
- An alternative: Computer-Assisted Clustering
  - Easy in theory: list all clusterings; choose the best
  - Impossible in practice: Too hard for us mere humans!
  - An organized list will make the search possible

# Switch from Fully Automated to Computer Assisted

- Fully Automated Clustering may succeed sometimes, but fails in general: too hard to understand when each model applies
- An alternative: Computer-Assisted Clustering
  - Easy in theory: list all clusterings; choose the best
  - Impossible in practice: Too hard for us mere humans!
  - An organized list will make the search possible
  - Insight: Many clusterings are perceptually identical

# Switch from Fully Automated to Computer Assisted

- Fully Automated Clustering may succeed sometimes, but fails in general: too hard to understand when each model applies
- An alternative: Computer-Assisted Clustering
  - Easy in theory: list all clusterings; choose the best
  - Impossible in practice: Too hard for us mere humans!
  - An organized list will make the search possible
  - Insight: Many clusterings are perceptually identical
  - E.g.,: consider two clusterings that differ only because one document (of 10,000) moves from category 5 to 6

# Switch from Fully Automated to Computer Assisted

- Fully Automated Clustering may succeed sometimes, but fails in general: too hard to understand when each model applies
- An alternative: Computer-Assisted Clustering
  - Easy in theory: list all clusterings; choose the best
  - Impossible in practice: Too hard for us mere humans!
  - An organized list will make the search possible
  - Insight: Many clusterings are perceptually identical
  - E.g.,: consider two clusterings that differ only because one document (of 10,000) moves from category 5 to 6
- Question: How to organize clusterings so humans can understand?

# Our Idea: Meaning Through Geography

Set of clusterings

# Our Idea: Meaning Through Geography

Set of clusterings $\approx$

A list of unconnected addresses

# Our Idea: Meaning Through Geography

Set of clusterings $\approx$

A list of unconnected addresses

Set of clusterings ≈
A list of unconnected addresses



⤳ We develop a (conceptual) geography of clusterings

# A New Strategy

Make it easy to choose best clustering from millions of choices

# A New Strategy

**Make it easy to choose best clustering from millions of choices**

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

# A New Strategy

Make it easy to choose best clustering from millions of choices

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

# A New Strategy
## Make it easy to choose best clustering from millions of choices

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

# A New Strategy
### Make it easy to choose best clustering from millions of choices

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

# A New Strategy
### Make it easy to choose best clustering from millions of choices

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

# A New Strategy
## Make it easy to choose best clustering from millions of choices

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

# A New Strategy
## Make it easy to choose best clustering from millions of choices

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

# A New Strategy
## Make it easy to choose best clustering from millions of choices

1. Code text as numbers (in one *or more* of several ways)
2. Apply all clustering methods we can find to the data — each representing different (unstated) substantive assumptions ($<$15 mins)
3. (Too much for a person to understand, but organization will help)
4. Develop an application-independent distance metric between clusterings, a metric space of clusterings, and a 2-D projection
5. "Local cluster ensemble" creates a new clustering at any point, based on weighted average of nearby clusterings
6. A new animated visualization to explore the space of clusterings (smoothly morphing from one into others)
7. ⇝ Millions of clusterings, easily comprehended
8. (Or, our new strategy: represent the entire bell space directly; no need to examine document contents)

# Many Thousands of Clusterings, Sorted & Organized
You choose one (or more), based on insight, discovery, useful information,...

# Software Screenshot

# Application-Independent Distance Metric: Axioms

- Metric based on 3 assumptions

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

- 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)

- 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(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 ⇒ 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(num clusters)

- ↝ Only <u>one</u> 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 ⇒ 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(num clusters)

- ⤳ 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:

- 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

- 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 ⇒ RA coders
    - Informative discoveries ⇒ 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

- What Are Humans Good For?

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

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

- **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**

- **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

- **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

- **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

- **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



(Our Method) – (Human Coders)

# Evaluation 1: Cluster Quality



Lautenberg Press Releases

(Our Method) – (Human Coders)

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



(Our Method) – (Human Coders)

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

- 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":

<u>Our Method 1</u> → vMF 1 → vMF 2 → <u>Our Method 2</u> → K-Means 1 → 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":

<u>Our Method 1</u> → vMF 1 → vMF 2 → <u>Our Method 2</u> → K-Means 1 → K-Means 2

"Genetic testing":

<u>Our Method 1</u> → {<u>Our Method 2</u>, K-Means 1, K-means 2} → Dir Proc. 1 → Dir Proc. 2

- David Mayhew's (1974) famous typology

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

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

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

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

- 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

# Example Discovery



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

# 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:

Space between methods:

Space between methods:
local cluster ensemble

# Example Discovery

# Example Discovery



Found a region with particularly insightful clusterings

# Example Discovery



Mixture:

Mixture:

0.39 Hclust-Canberra-McQuitty

Mixture:

0.39 Hclust-Canberra-McQuitty

0.30 Spectral clustering
Random Walk
(Metrics 1-6)

Mixture:

0.39 Hclust-Canberra-McQuitty

0.30 Spectral clustering
Random Walk
(Metrics 1-6)

0.13 Hclust-Correlation-Ward

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

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

Mixture:

0.39 Hclust-Canberra-McQuitty

0.30 Spectral clustering
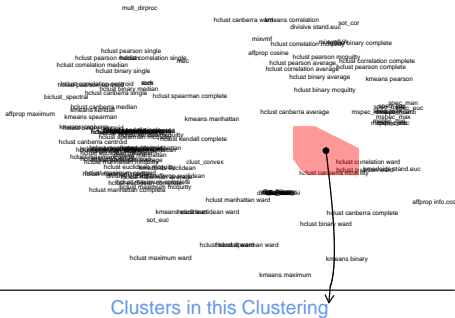Random Walk
(Metrics 1-6)

0.13 Hclust-Correlation-Ward

0.09 Hclust-Pearson-Ward

0.05 Kmediods-Cosine

0.04 Spectral clustering
Symmetric
(Metrics 1-6)

# Example Discovery



Clusters in this Clustering

Clusters in this Clustering

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"

Mayhew

Clusters in this Clustering

Credit Claiming
Pork

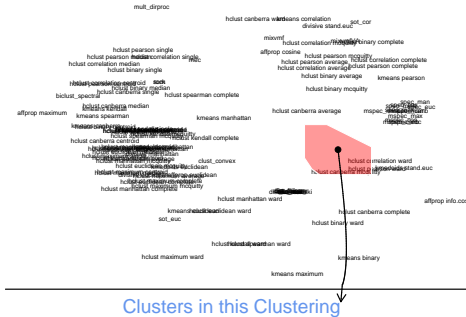Mayhew  Credit Claiming
Legislation

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
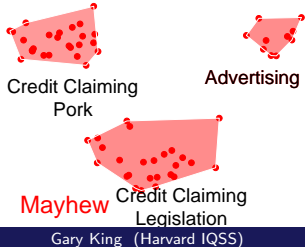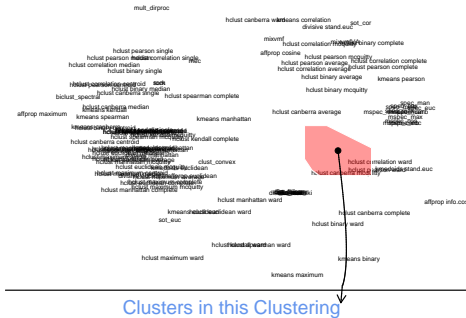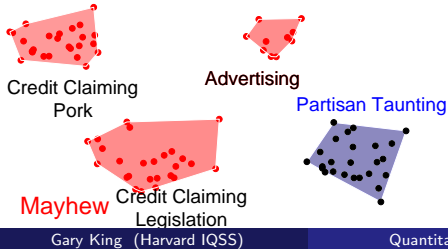


Advertising:
"Senate Adopts Lautenberg/Menendez Resolution Honoring Spelling Bee Champion from New Jersey"

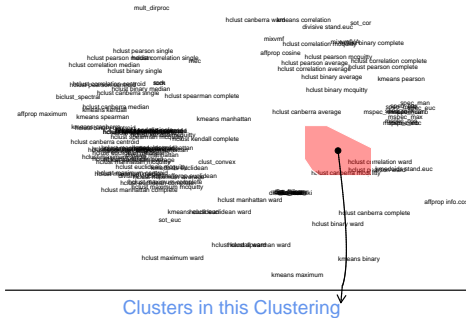# Example Discovery: Partisan Taunting



Clusters in this Clustering

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

Credit Claiming Pork

Advertising

Partisan Taunting

Mayhew

Credit Claiming Legislation

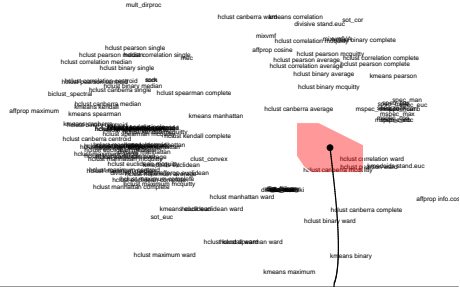# Example Discovery: Partisan Taunting



Clusters in this Clustering

Credit Claiming Pork

Advertising

Partisan Taunting

Mayhew

Credit Claiming Legislation

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

# Example Discovery: Partisan Taunting



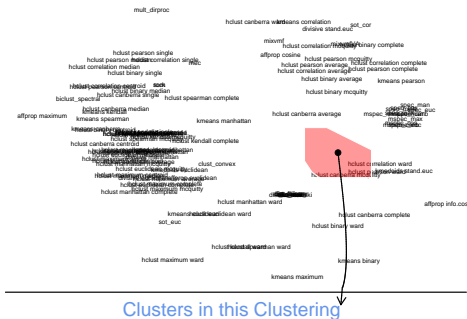Clusters in this Clustering

Credit Claiming Pork

Advertising

Partisan Taunting

Mayhew

Credit Claiming Legislation

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

Taunting ruins deliberation

# 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]

## 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]
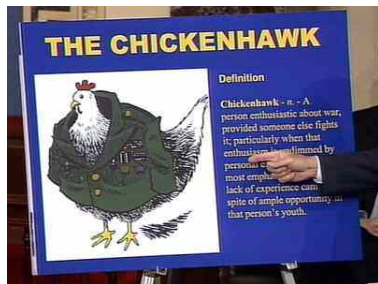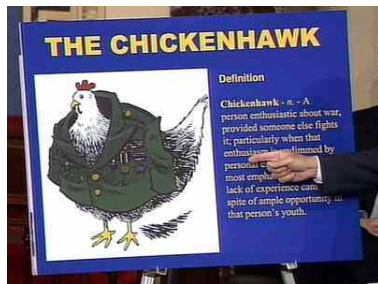
## Taunting ruins deliberation



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]

- Discovered using 200 press releases; 1 senator.

# Out of Sample Confirmation of Partisan Taunting

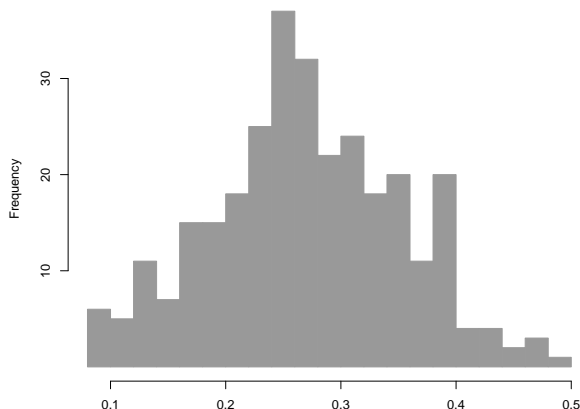- Discovered using 200 press releases; 1 senator.
- Confirmed using 64,033 press releases; 301 senator-years.

# 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 <span style="color:red">proportion of press releases</span> 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
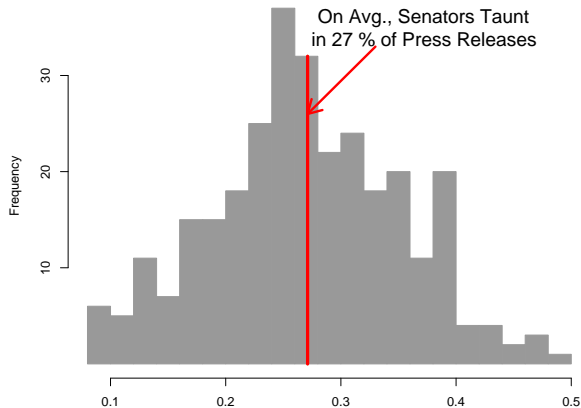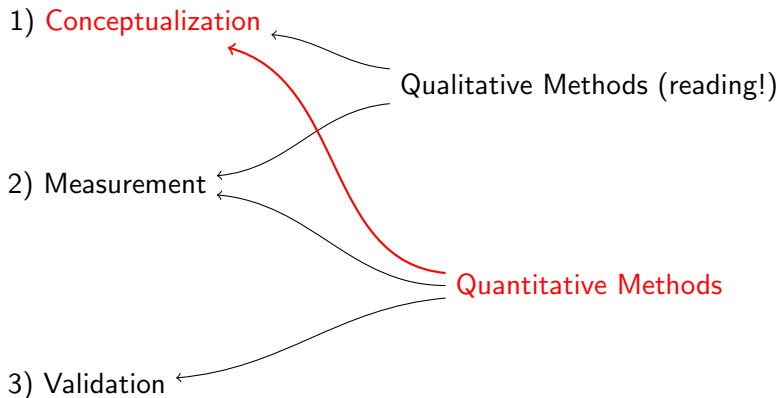
# Out of Sample Confirmation of Partisan Taunting

- 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



On Avg., Senators Taunt in 27 % of Press Releases

# Quantitative Methods for Qualitative Conceptualization



1) Conceptualization

Qualitative Methods (reading!)

2) Measurement
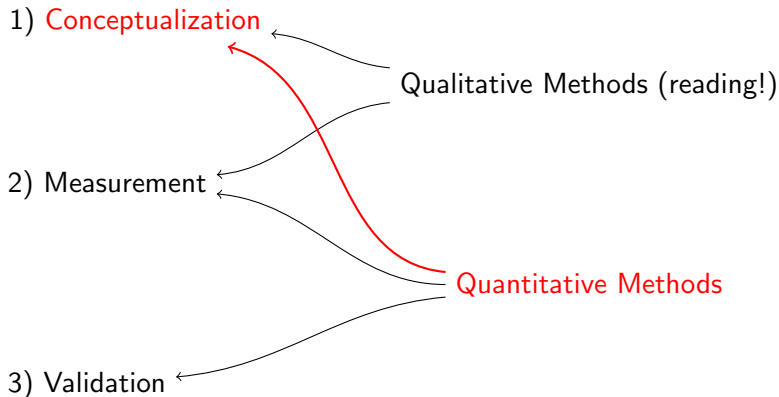
Quantitative Methods

3) Validation

Quantitative methods for conceptualization and discovery

# Quantitative Methods for Qualitative Conceptualization



1) Conceptualization

Qualitative Methods (reading!)

2) Measurement

Quantitative Methods

3) Validation

Quantitative methods for conceptualization and discovery
- Few formal methods designed explicitly for conceptualization

# Quantitative Methods for Qualitative Conceptualization
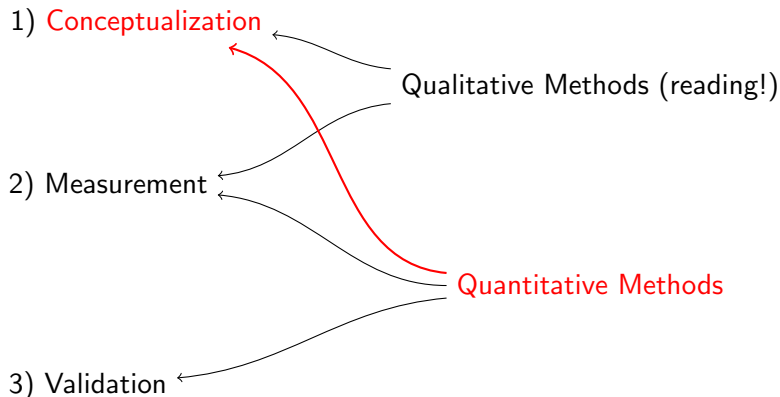
1) Conceptualization

Qualitative Methods (reading!)

2) Measurement

Quantitative Methods

3) Validation

Quantitative methods for conceptualization and discovery

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

# Quantitative Methods for Qualitative Conceptualization



1) Conceptualization

Qualitative Methods (reading!)

2) Measurement

Quantitative Methods

3) Validation

Quantitative methods for conceptualization and discovery

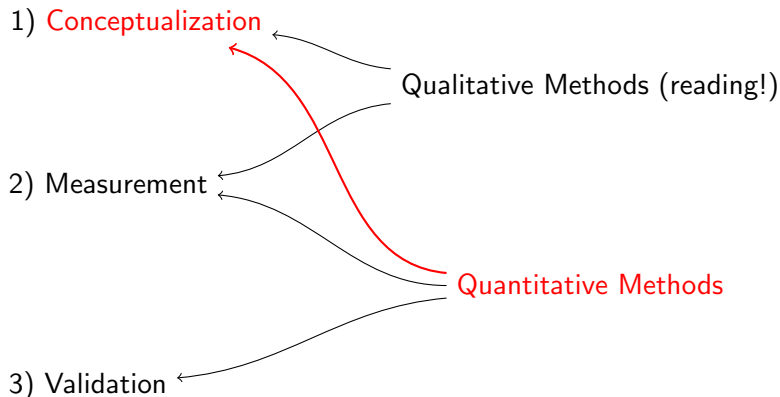- Few formal methods designed explicitly for conceptualization
- Belittled: "Tom Swift and His Electric Factor Analysis Machine" (Armstrong 1967)
- Evaluation methods measure progress in discovery

# For more information



http://GKing.Harvard.edu