

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

Gary King

Institute for Quantitative Social Science  
Harvard University

Talk on 2/5/2010 for All-hands IQSS Staff Meeting

Joint work with Justin Grimmer, Harvard University

# Some context for related technology

- <http://ow.ly/14hDU>
- <http://ow.ly/14h36>

# A Method for Conceptualization

- Systematic method for computer-assisted conceptualization from text

# A Method for Conceptualization

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

# A Method for Conceptualization

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

# A Method for Conceptualization

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

# Why Johnny Can't Classify (Optimally)

# Why Johnny Can't Classify (Optimally)

- Clustering seems easy; its not!

# Why Johnny Can't Classify (Optimally)

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

# Why Johnny Can't Classify (Optimally)

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

# Why Johnny Can't Classify (Optimally)

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

# Why Johnny Can't Classify (Optimally)

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

# Why Johnny Can't Classify (Optimally)

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

# Why Johnny Can't Classify (Optimally)

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

# Why Johnny Can't Classify (Optimally)

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

# Why Johnny Can't Classify (Optimally)

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

# Why HAL Can't Classify Either

# Why HAL Can't Classify Either

- Large quantitative literature on **cluster analysis**

# Why HAL Can't Classify Either

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

# Why HAL Can't Classify Either

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

# Why HAL Can't Classify Either

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

# Why HAL Can't Classify Either

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

# Why HAL Can't Classify Either

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

# Why HAL Can't Classify Either

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

# Why HAL Can't Classify Either

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

# Why HAL Can't Classify Either

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

# Why HAL Can't Classify Either

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

# Our Idea: Meaning Through Geography

Set of clusterings

# Our Idea: Meaning Through Geography

Set of clusterings  $\approx$

A list of unconnected addresses

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

	195	Car	C
<b>Cartage New England Inc</b> 28 Allen Ln Ipswich 01938..... 978 356-9960	<b>Carter F</b> 34 Hibiscus Bldg 02133..... 617 327-1105	<b>Carter Nella E</b> 323 Main St 02115..... 617 267-6483	
<b>Cartagena Lydia</b> 20 Sweet Box 02131..... 617 323-7639	<b>Faye &amp; Ricky</b> 20 Columbia Ave Box 02136..... 617 437-7331	<b>Nicholas S F</b> 115 Randolph Ave Mill 02186..... 617 698-5307	
<b>Cartagena Avish</b> F Pleasant Box 02139..... 617 442-9780	<b>Francis S</b> 134 Yankov W Ave 02132..... 617 323-6781	<b>Nick 21 Farwell Box 02114..... 617 267-5222</b>	
<b>B Had 02134</b> ..... 617 361-5253	<b>Franklin &amp; Anne</b> 201 Mt Auburn Cam 02138..... 617 354-0798	<b>Nick &amp; Debbi</b> 196 Morris Rd Newton 02459..... 617 527-0480	
<b>17 566-1282</b> <b>Jessica</b> 50 Decatur Cha 02129..... 617 241-0152	<b>Fred 41 Howland Elm 02136..... 617 524-3078</b>	<b>Nicole..... 617 698-0713</b>	
<b>17 364-5188</b> <b>Luzmila</b> 124 Harvard Cam 02138..... 617 491-5621	<b>Fred 16 Howland Ave Mill 02136..... 617 698-1343</b>	<b>Norman G</b> 38 Chickareed Dr 02125..... 617 822-1201	
<b>361-0380</b> <b>Melvin</b> 503 Green Cam 02129..... 617 576-1061	<b>G &amp; B</b> 8 Verden Dcr 02134..... 617 436-8906	<b>P 40 Cranford Pl Box 02135..... 617 437-4754</b>	
<b>17 566-4548</b> <b>Carl Nicholas</b> 18 Appleton Boston 02114..... 617 695-6996	<b>G T 27 Franklin Ave Sun 02145..... 617 623-7121</b>	<b>P E 501 E South S Box 02137..... 617 268-8213</b>	
<b>17 628-8248</b> <b>Carlton</b> 0 4 Bradford Box 02133..... 617 338-9219	<b>George 125 Madison Box 02134..... 617 367-9548</b>	<b>P 81 16 Hutchings Box 02131..... 617 427-9170</b>	
<b>17 445-5116</b> <b>Thomas &amp; Kathleen</b> 17 Franklin St Mill 02136..... 617 698-6163	<b>Carter Hillside Assoc</b> 107 S Street Box 02111..... 617 456-1689	<b>Paul &amp; Constance</b> 114 Franklin St W Box 02131..... 617 325-2036	
<b>17 822-2962</b> <b>Carter A</b> 50 Thompson Ln Mill 02136..... 617 696-6919	<b>Carter Harry F</b> 30 Bayview Rd W Box 02132..... 617 325-5465	<b>Paul M 501 E South S S Box 02137..... 617 268-4546</b>	
<b>17 427-5712</b> <b>A Nelson</b> A 201 Beulah Ave Cambridge 02238..... 617 492-4174	<b>Carter Hide Co Inc</b> 161 Huntington St 02148..... 617 542-7987	<b>Paul M 27 Crown Bk 02139..... 617 787-2115</b>	
<b>17 569-2698</b> <b>A 21 Beulah Way Rosbury 02139..... 617 442-1219</b>	<b>Carter Hilary 41 Harvey Cam 02148..... 617 876-2750</b>	<b>Prangman 02102..... Wellesley Tpk 781.235-0488</b>	
<b>17 667-5190</b> <b>A M 255 Massachusetts Ave 02115..... 617 266-7153</b>	<b>Horace</b> 301 Walnut Ave Rosbury 02139..... 617 442-5307	<b>Carter Prudence</b> 40 Franklin Waterbury 02172..... 617 393-3782	
<b>17 569-1412</b> <b>Adams 361 Centre St Mill 02136..... 617 698-7074</b>	<b>Howard Jr 28 Nona Drive Box 02118..... 617 445-5532</b>	<b>Prudence</b> 40 Franklin Waterbury 02172..... 617 926-7063	
<b>17 338-9110</b> <b>Alice 108 Elmwood Box 02134..... 617 423-0193</b>	<b>J Dan..... 617 354-2658</b>	<b>Roginald</b> 106 Brookview Dorchester 02122..... 617 541-2843	
<b>17 825-1993</b> <b>Alice 40 Market Cambridge 02238..... 617 945-2711</b>	<b>J 23 Chatham Box 02446..... 617 232-7990</b>	<b>Renee &amp; Andrew</b> 30 Walnut Box 02138..... 617 720-3765	
<b>17 296-1593</b> <b>Andrew F 42 West St Sun 02133..... 617 625-7623</b>	<b>J 538 Harvard Box 02446..... 617 730-9483</b>	<b>Carter Rice David</b> 3000 Weston Publishing 163 Main Wilmington 01887	
<b>17 670-2078</b> <b>Carter Anne MD</b> 1161 Beacon Box 02446..... 617 739-1022	<b>Jacques J 775 The Pines West Rosbury 02132..... 617 323-5374</b>	<b>Ted Free-Dial '3' &amp; Thru..... 800 638-1671</b>	
<b>17 621-9001</b> <b>Carter J M</b> 310 Columbia Rd S Box 02137..... 617 464-1040	<b>Carter J D</b> 40 Newmarket Pl Box 02146..... 617 876-8841	<b>Ted Free-Dial '3' &amp; Thru..... 800 616-7447</b>	
<b>17 296-4725</b> <b>B E 108 Graduate Ave Mill 02136..... 617 296-6911</b>	<b>Carter J Venal Co</b> 40 Newmarket Pl Box 02146..... 617 442-1775	<b>Ted Free-Dial '3' &amp; Thru..... 800 648-7447</b>	
<b>17 542-1521</b> <b>Carter Barbara L MD</b> Turks-New England Medical Center Box 02111	<b>Carter James</b> 157 Cambridge St Cam 02136..... 617 492-1214	<b>Call..... 978 988-7447</b>	
<b>17 364-5232</b> <b>Carter Broadcasting Co</b> 30 Park Pl Box 02134..... 617 423-4368	<b>James 412 Foster Ave Rosbury 02138..... 617 739-2193</b>	<b>Ingalls Centre 163 Main Wilmington 01887..... 800 638-1673</b>	
<b>17 739-2662</b> <b>Bernard J</b> 122 Goodville E Box 02136..... 617 567-3430	<b>James 130 Good Star Rd Cambridge 02141..... 617 876-8841</b>	<b>Richard A</b> 2079 Lawrence Ave Brighton 02215..... 617 982-0836	
<b>17 879-0030</b> <b>Bibbiah 25 Midway Dcr 02124..... 617 298-8713</b>	<b>Jane L 34 Rosbury Rd Mill 02134..... 617 361-0773</b>	<b>Richard A MD</b> 47 Mt Vernon Box 02106..... 617 566-7293	
<b>17 541-5249</b> <b>Bill 26 Mt Vernon Box 02106..... 617 367-9931</b>	<b>Janice 14 Adams Rd Newton 02458..... 617 564-0435</b>	<b>Richard R K</b> 130 Canterbury Ln 02126..... 617 267-0710	
<b>17 879-0030</b> <b>Carter &amp; Business Consultants Inc</b> 30 Park Pl Box 02134..... 617 423-4310	<b>Jacques 41 Warren Ave Box 02134..... 617 424-5994</b>	<b>Richard R K</b> 23 Mower S Box 02137..... 617 268-0448	
<b>17 436-1511</b> <b>Carter C 200 Commonwealth Ave 02135..... 617 782-2118</b>	<b>John 107 Summer Box 02125..... 617 423-4334</b>	<b>Roger 130 St Bourne Box 02131..... 617 424-6148</b>	
<b>17 569-4119</b> <b>C 218 Harvard Ave East Boston 02268..... 617 569-1545</b>	<b>John 40 Howland Elm 02136..... 617 282-1235</b>	<b>Roy 41 Concord Cam 02138..... 617 491-6115</b>	
<b>800 569-8782</b> <b>C 109 Harvard Cam 02138..... 617 491-8522</b>	<b>June O 129 A Summit Ave 02133..... 617 734-6109</b>	<b>Royce 18 Sundry Cha 02129..... 617 241-0418</b>	
	<b>J 29 Howland Elm 02136..... 617 282-1232</b>		
	<b>K 17 Concord Dorchester 02127..... 617 282-1293</b>		





# 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)

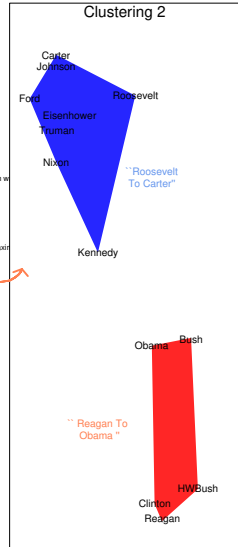
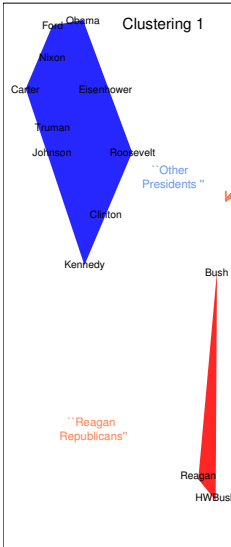
# A New Strategy

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)

# Many Thousands of Clusterings, Sorted & Organized

You choose one (or more), based on insight, discovery, useful information, . . .



# Application-Independent Distance Metric: Axioms

# Application-Independent Distance Metric: Axioms

- Metric based on 3 assumptions

# Application-Independent Distance Metric: Axioms

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

# Application-Independent Distance Metric: Axioms

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

# Application-Independent Distance Metric: Axioms

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

# Application-Independent Distance Metric: Axioms

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

# Application-Independent Distance Metric: Axioms

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

# Evaluating the Performance of Our Method

# Evaluating the Performance of Our Method

- Goals:

# Evaluating the Performance of Our Method

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

# Evaluating the Performance of Our Method

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

# Evaluating the Performance of Our Method

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

# Evaluating the Performance of Our Method

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

# Evaluating the Performance of Our Method

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

# Evaluating the Performance of Our Method

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

# Evaluating the Performance of Our Method

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

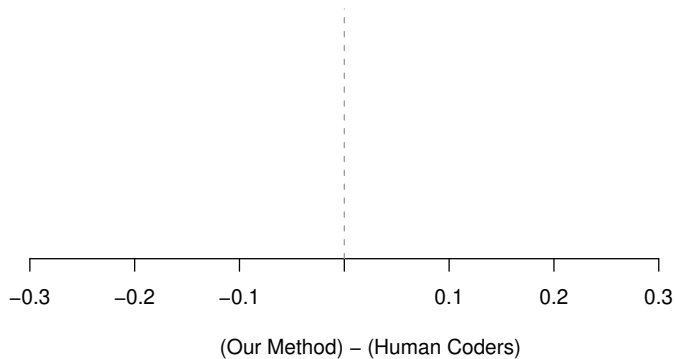
# Evaluation 1: Cluster Quality

- **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**

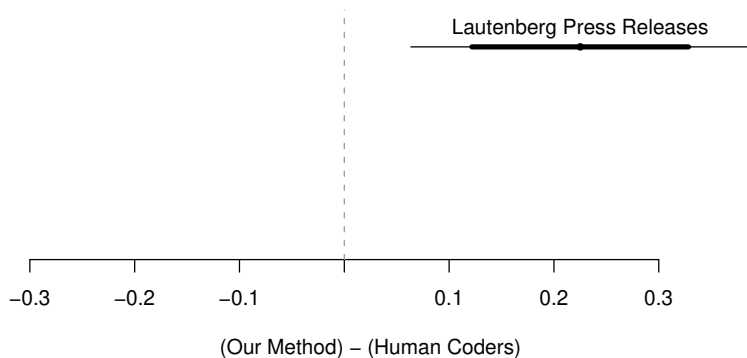
# Evaluation 1: Cluster Quality

- **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

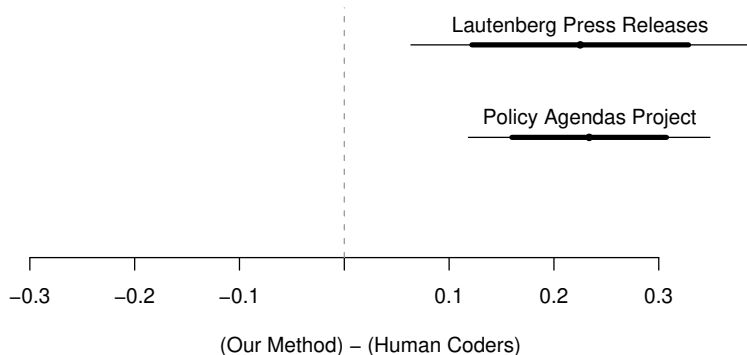


# Evaluation 1: Cluster Quality



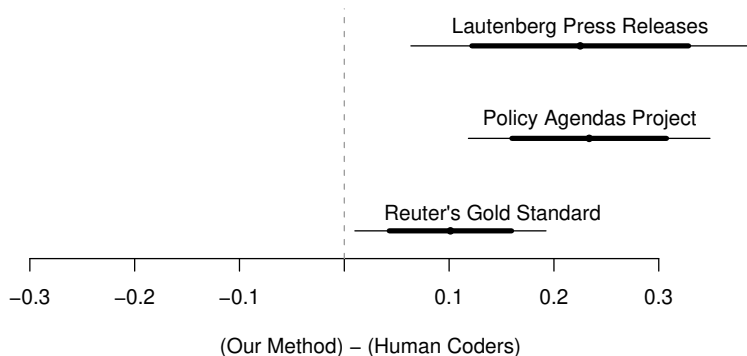
Lautenberg: 200 Senate Press Releases (appropriations, economy, education, tax, veterans, ...)

# 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

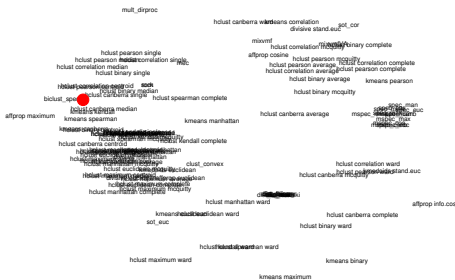








# Example Discovery



Space between methods:







# Example Discovery



Mixture:













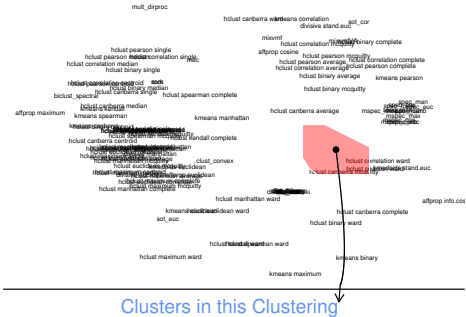




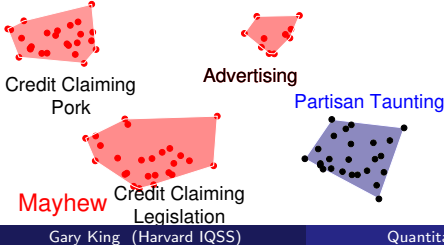




# Example Discovery: Partisan Taunting



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

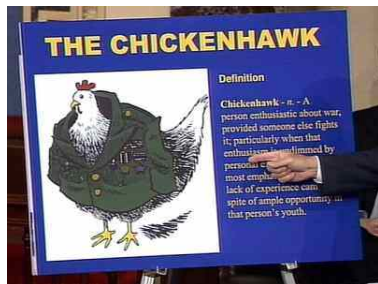








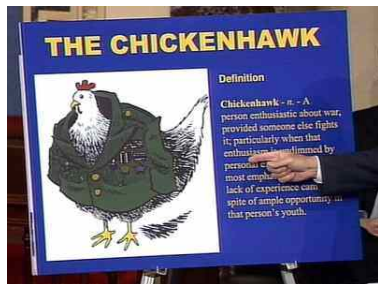
## Taunting ruins deliberation



Sen. Lautenberg  
on Senate Floor  
4/29/04

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

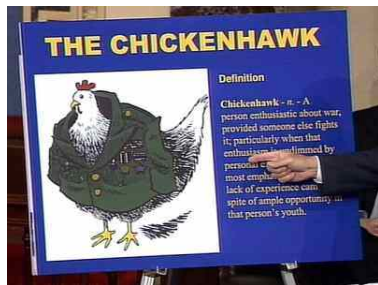
## Taunting ruins deliberation



Sen. Lautenberg  
on Senate Floor  
4/29/04

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

## Taunting ruins deliberation



Sen. Lautenberg  
on Senate Floor  
4/29/04

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

# Out of Sample Confirmation of Partisan Taunting

- Discovered using 200 press releases; 1 senator.

# Out of Sample Confirmation of Partisan Taunting

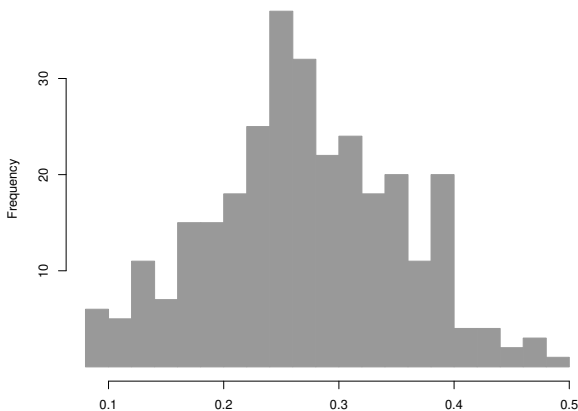
- Discovered using 200 press releases; 1 senator.
- Confirmed using 64,033 press releases; 301 senator-years.

# Out of Sample Confirmation of Partisan Taunting

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

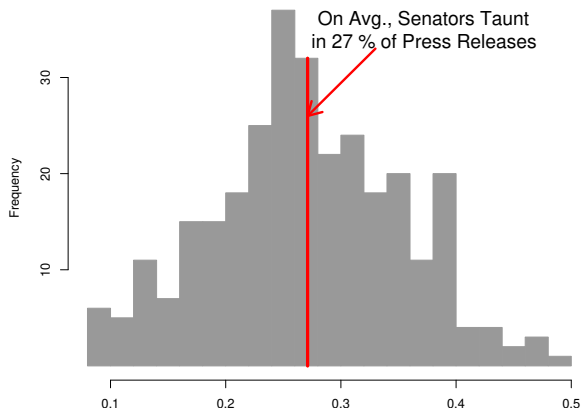
# Out of Sample Confirmation of Partisan Taunting

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

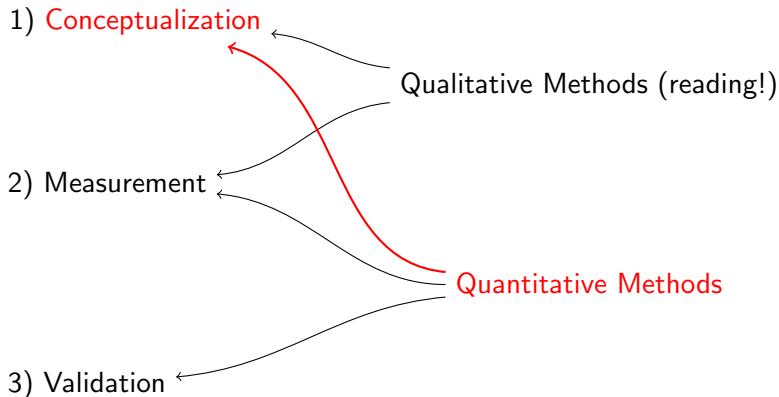


# Out of Sample Confirmation of Partisan Taunting

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

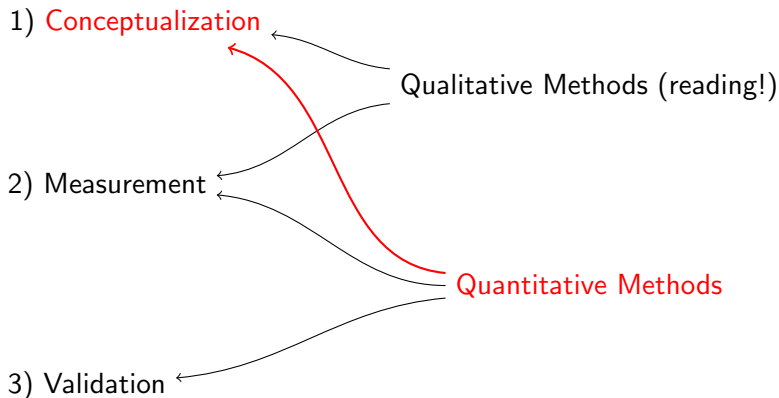


# Advancing the Objective of Discovery



Quantitative methods for conceptualization: aiding **discovery**

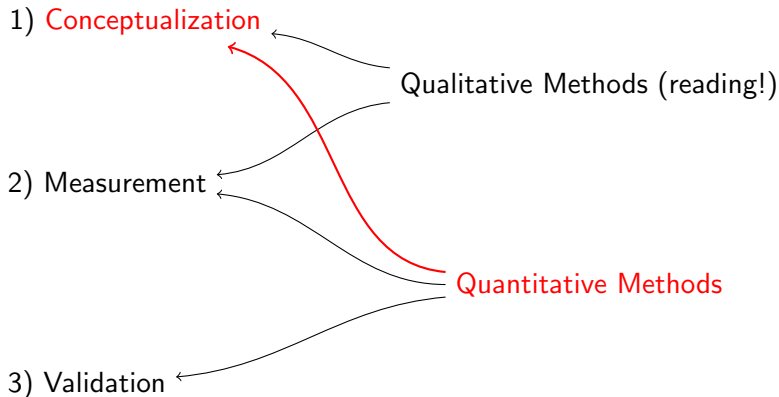
# Advancing the Objective of Discovery



Quantitative methods for conceptualization: aiding **discovery**

- Few formal methods designed explicitly for conceptualization

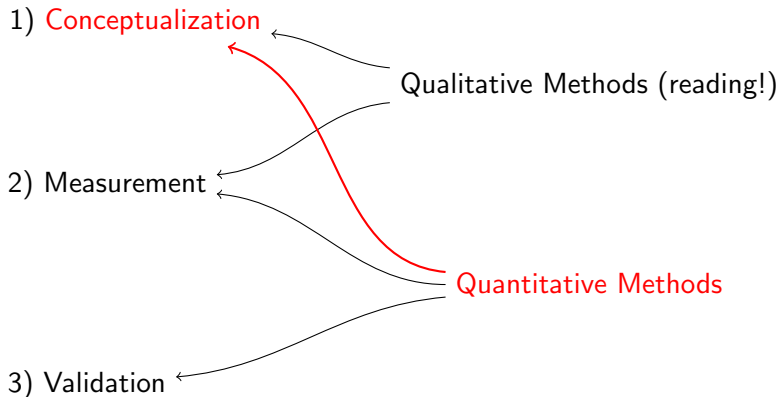
# Advancing the Objective of Discovery



Quantitative methods for conceptualization: aiding **discovery**

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

# Advancing the Objective of Discovery



Quantitative methods for conceptualization: aiding **discovery**

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

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

<http://GKing.Harvard.edu>