

# The Devil You Know: Threat Construction in the United States

1

**DATA SOURCE:  
INTERNET ARCHIVE'S .GOV COLLECTION**

**E-SCIENCE DATA INCUBATOR PARTICIPANT:  
EMILY KALAH GADE**

**PROJECT LEADS:  
DR. JOHN WILKERSON  
DR. MAGDALENA BALAZINSKA**

# Project Background: Scope

2

- Political Science Side - NSF Grant re: Poli-Infomatics
  - Explore the wealth of textual resources available today
    - ✦ Measuring ideas in text... Using computers to do it
- Computer Science Side – Text in Myria

# Project Background: Data Source

3

- Internet Archive .gov collection
  - 90 terabytes of data, 1.1 billion captures
  - Hosted by Altiscale on a Hadoop Cluster
- Incubator project:
  - Figure out what is “in there”
  - Make it useful for social science research
    - ✦ Focus: Language about threats – terrorism, climate change, and banking crisis

# Incubator Goals

4

- Access the cluster, run scripts, get output
- Reliability and validity
- Results and write up
- Make it reproducible and accessible to other social scientists

# Initial Hurdles

5

- Emily's coding abilities
  - General big data issues
  - Nit picky issues: unicode, parsers, carriage returns
  - Best format for examining data
  - Reliability and validity
- 
- Reproducibility for anyone else (without the incubator!)

## 6

12/5/14

# What word count data looks like

7

2009	5	dot\.	desertification	1	
2009	5	hhs\.	(fresh\swater)	21	
2009	5	state\.	desertification	1409	
2009	5	state\.	(climate\schange)		28204
2009	5	energy\.	pollution	5812	
2009	5	energy\.	anthropocene	1	
2009	5	\.house\.	total	655546606	
2009	5	defense\.	pollution	93	
2009	6	ed\.	pollution	732	
2009	6	va\.	(global\swarming)	3	
2009	6	dod\.	(food\ssecurity)	3	
2009	6	dod\.	(natural\sdisaster)		274
2009	6	dol\.	(forest\sconservation)	9	
2009	6	dot\.	(food\ssecurity)	25	
2009	6	dot\.	(greenhouse\sgas)	941	
2009	6	hhs\.	pollution	1151	
2009	6	state\.	(food\ssecurity)		4959
2009	6	state\.	(natural\sdisaster)		9621
2009	6	\.senate\.	anthropogenic	526	
2009	6	\.senate\.	(ocean\sacidification)	48	
2009	6	whitehouse\.	(global\swarming)		281
2009	7	other	(greenhouse\sgas)	138771	
2009	7	va\.	(climate\schange)	16	
2009	7	dod\.	(fresh\swater)	55	
2009	7	state\.	(fresh\swater)	838	
2009	7	energy\.	(forest\sconservation)	2	
2009	7	\.house\.	(security\sof\sfood)		147
2009	7	\.senate\.	total	545452207	
2009	7	whitehouse\.	desertification	1	
2009	7	whitehouse\.	(climate\schange)		2157
2009	8	va\.	total	100787717	
2009	8	va\.	(natural\sdisaster)	492	

# Data Quality

8

- Complete Crawls Taken from Nov to Jan around elections from 2004 to present
  - This becomes Gold Standard
- Early data crawled irregularly
  - Early data sensitive to terms



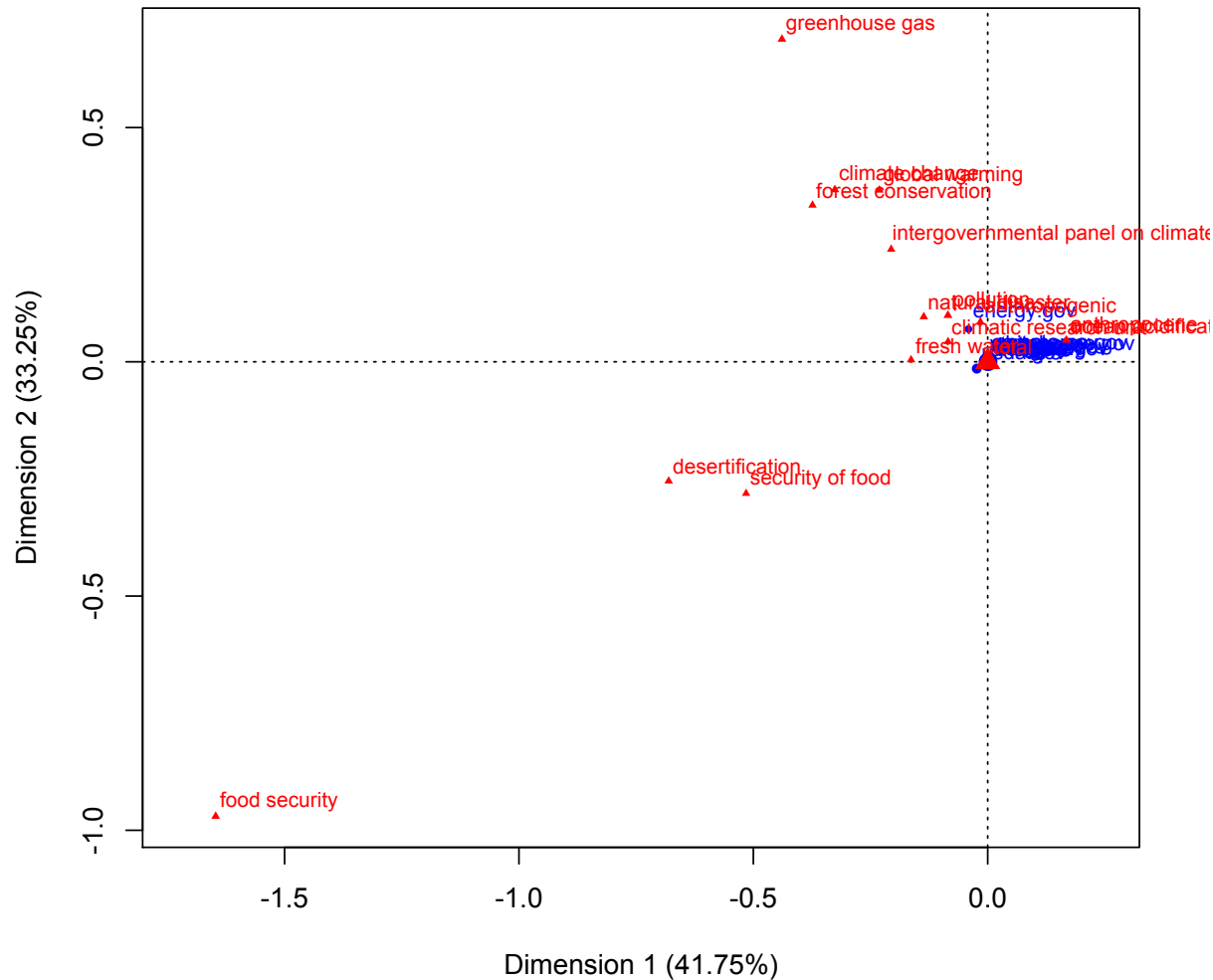
# Preliminary Results

9

- **Correspondence Analysis**
  - Angles between topics and agency indicate more emphasis of that topic by that agency (smaller angles = more emphasis)
  - Interagency distances and inter topic distances are interpretable (but not topic-agency distances)

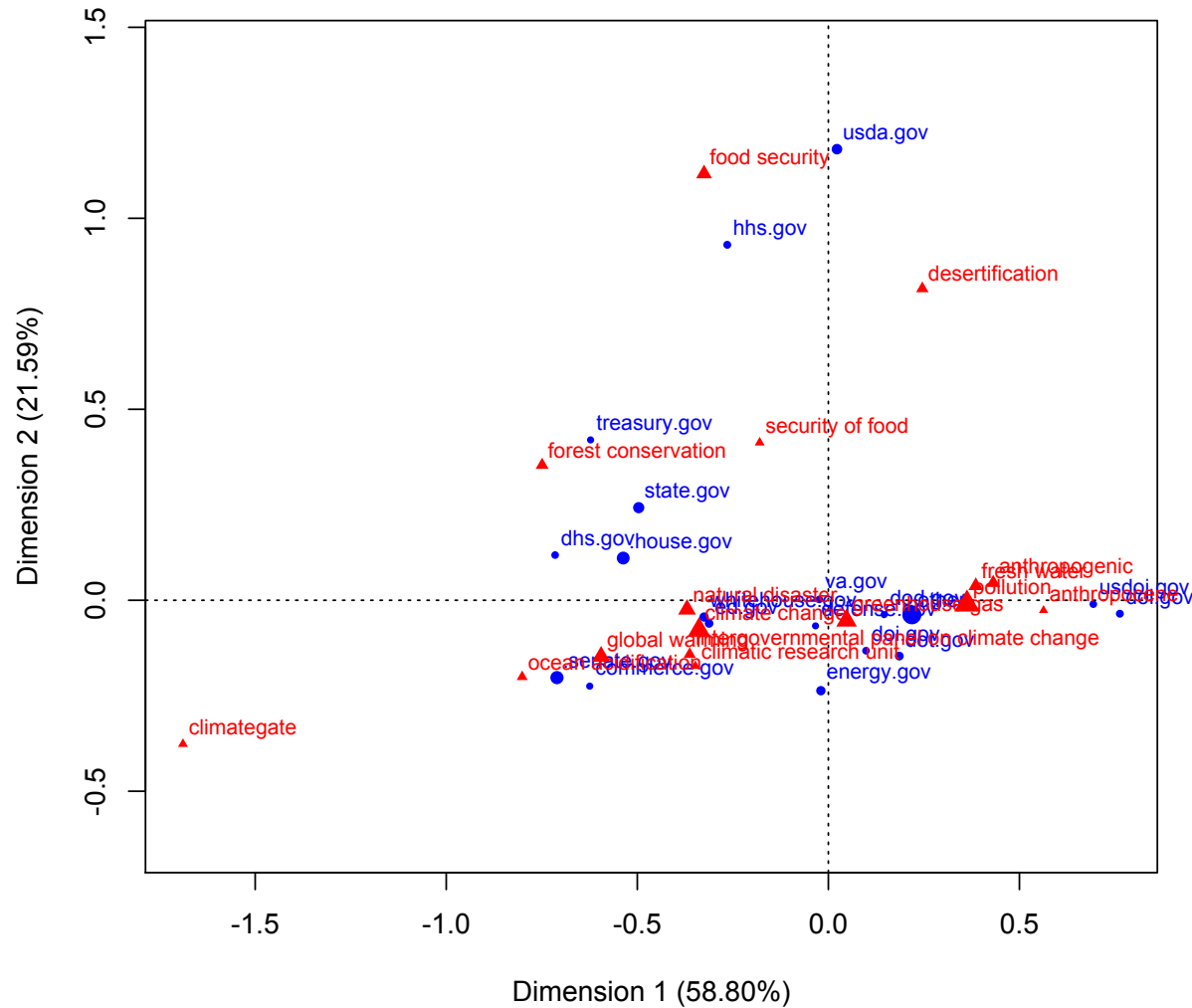
# 2004 Agency Relative Emphasis Biplot

10



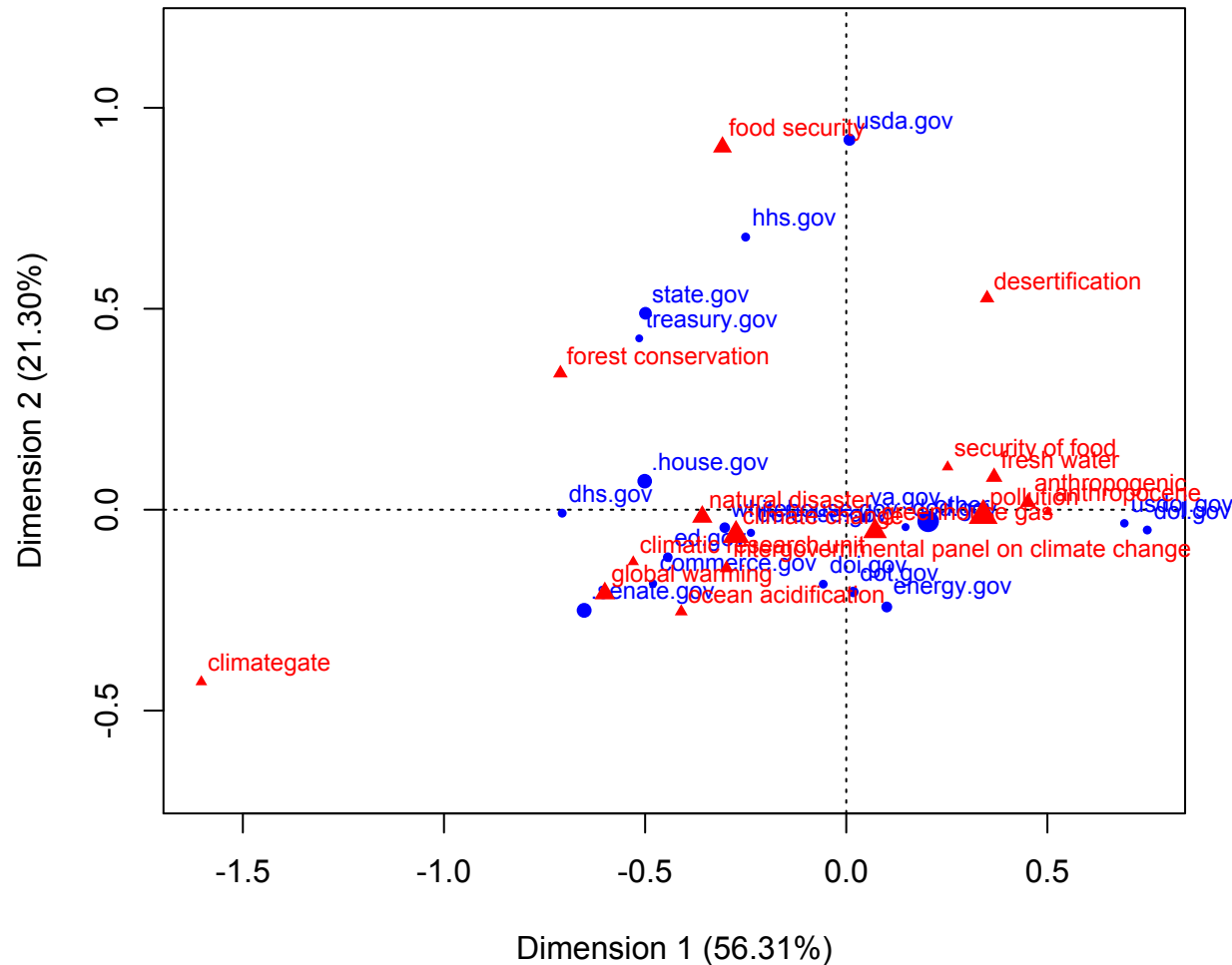
# 2008 Agency Relative Emphasis Biplot

11



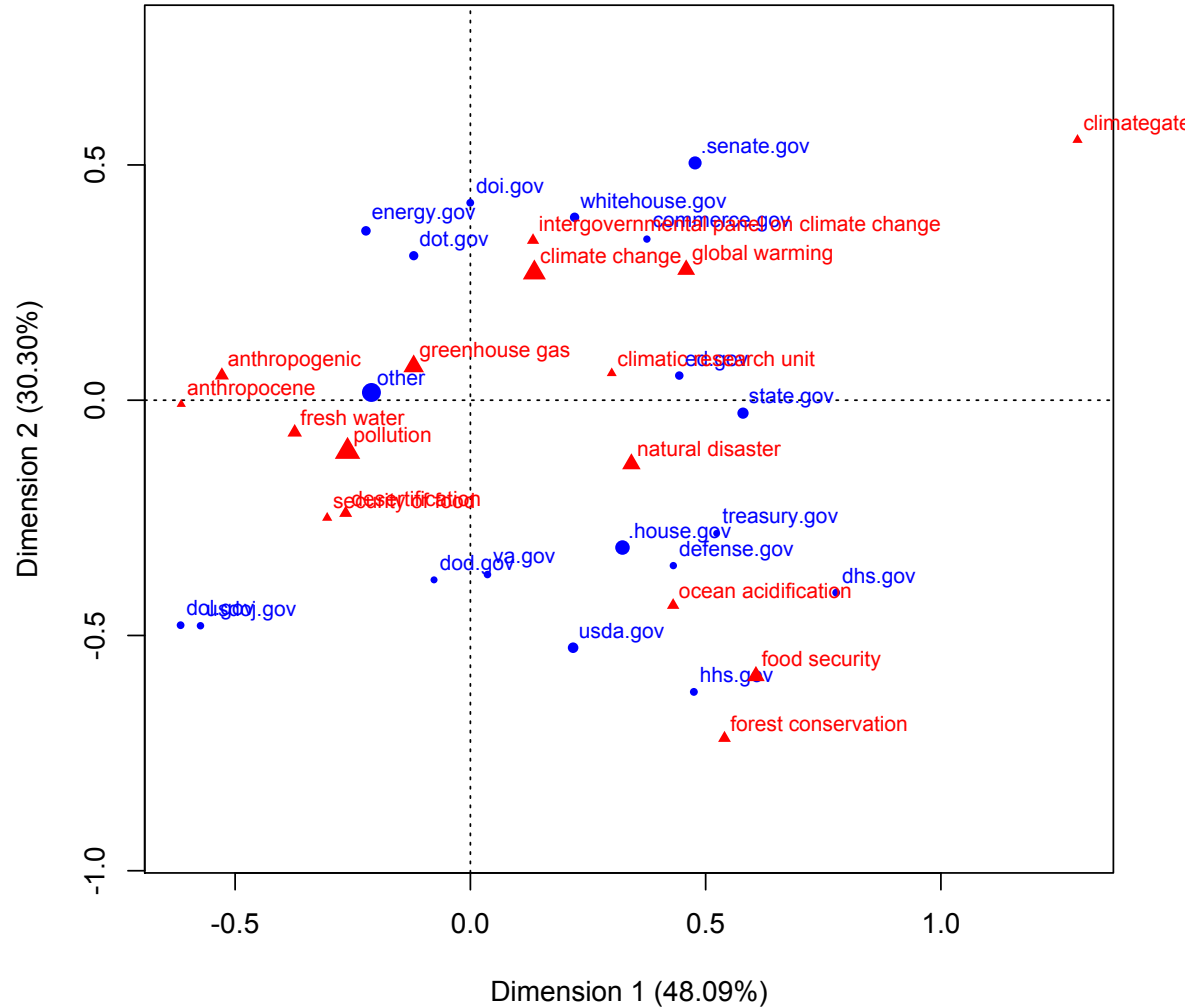
# 2012 Agency Relative Emphasis Biplot

12



# 2013 Agency Relative Emphasis Biplot

13



# Creating Institutional Memory

14

- Taking data to PoliInfomatics conference
- Publishing code on Github, publishing data on PoliInfomatics
- Contact with Internet Archive and two other groups known to be exploring this data
  - One at Harvard, one in Germany
- Article re: .gov collection as a data source

# Proposed Publications

15

- Methodological Paper re: .gov data as a data source for social science research
  - Including preliminary results and scripts
- The Devil You Know: Threat Construction in the United States
  - Tracking threat construction in the US; Agenda analysis

# Thank You!

16