Human Geography: Open Source Data, Methods and Forecasting

Philip Schrodt Political Science and International Center for the Study of Terrorism Pennsylvania State University schrodt@psu.edu http://eventdata.psu.edu

4 September 2012

What's the buzz? Why 2012 is not 2002



Political Data Sources





An Open-Source Application for Publishing, Citing and Discovering Research Data

> Universität Konstanz







Political Data Collection: Penn State

- Correlates of War Project (Bennett, Palmer)
- Legislative Text (Monroe)
- Event Data (Schrodt)
- World Religions (Finke, Henderson, Maoz)
- Association of Religion Data Archives (Finke)
- International Center for the Study of Terrorism (Horgan, Bloom)
- NSF Integrated Graduate Research and Training: Big Data in the Social Sciences (Monroe)

Computing Power

Control Data Corporation 3600 (ca.1965) 32 K (48-bit) RAM memory 1 processor ~1-million operations per second Output: line printer







Penn State High Performance Computing Facility 15 cluster computers 100 to 2000 2.66 Ghz processors in each cluster ~50 Gb RAM accessible to each processor 130 Tb disk space 4 interactive visualization rooms

Computing Power

Control Data Corporation 3600 (ca.1965) 32 K (48-bit) RAM memory 1 processor ~1-million operations per second Output: line printer







Penn State High Performance Computing Facility 15 cluster computers 100 to 2000 2.66 Ghz processors in each cluster ~50 Gb RAM accessible to each processor 130 Tb disk space 4 interactive visualization rooms



Motorola Razr 16 Gb RAM memory Dual processor ~500-million operations per sec 540 x 860 color display

Open Source Software



Event Data in 2012

- Combining the web, computing power and automated coding methods, we can now create instruments to automatically monitor the state of the political world in near real time.
- We are only beginning to develop the appropriate ways of systematically using this information

[In the slides which follow, the "Notes" format provides the source and sometimes additional information and/or commentary]

Statistics Packages :: Cars

Stata





R

SAS



SPSS



Statistics Packages :: Cars

Stata





R

MatLab



SAS



SPSS







• Open source



- Open source
- Widely used in all approaches to statistical analysis and pattern recognition
- CRAN library provides almost immediate access to new methods
 - events (event data aggregation)
 - tm (text miner)
- Robust scripting capabilities; easily interfaces with C/C++ when needed
- Skill set is widely available

Integrated open, real time data generation



Levant Event Data Set



Material Cooperation

Material Conflict

00-01

PSE to ISR , inverted scale

10-01

90-01



Time

GDELT Event Data

- Global coding, Jan-1979 to Jun-2012
- 180-million events based on open news sources
- CAMEO event, actor and sub-state agent coding
 - 15,000 verb phrase dictionary
 - 40,000+ political actors and agents
- Geolocated to city level
- Planned quarterly updates, backfit to 1900 or possibly 1800



GDELT: Coverage in 1% sample



ACLED: Actors across time





Figure 14: Violence against civilians by actor type, 1997 - 2012.

Figure 4: Count of discrete non-state actors, DR-Congo, Jan - Jul 2012.

ACLED: Geospatial



ACLED: Geospatial



European Media Monitor



Uppsala Georeferenced Event Dataset



The GED is the product of two and a half years of work at the **Department of Peace and Conflict Research, Uppsala University**. The UCDP GED contains conflict data disaggregated spatially and temporally down to the level of the individual incidents of violence. For more details please see the About UCDP link above.

Challenges to integrating models into decision-making

- Forecasting is hard (Tetlock)
- Probabilistic reasoning is hard (Kahneman, Taleb)
- Statistics is new compared to deterministic modeling and is still changing, even at very fundamental levels
 - Frequentist vs Bayesian approaches
 - New approaches made possible by computational advances
- USG has locked-in to a set of approaches which are 20 to 40 years old (oh, and don't work...)
 - The answers aren't simple, even if some colonel wants them to be simple
 - Our 20th century peer competitors were trained as political ideologues; our 21st century peer competitors are trained as engineers

IC Allocation of Resources by Methodology



IC Allocation of Resources by Methodology



Methods of Modeling

- Machine Learning
 - Very flexible methods of determining relationships
 - No [widely accepted] theory of error
- Classical ("frequentist") statistics
 - Objective is determining whether a variable has a non-zero effect: "significance tests"
 - Effective in experimental and other randomized settings, but generally useless for predictive models
- Bayesian statistics
 - Objective is determining a probability of the impact of a variable based on your prior knowledge and the new data
 - Corresponds to how most people actually think about data but has only become computationally feasible in the past twenty years

Statistical challenges

- Systematically dealing with measurement error and missing values rather than assuming "missing at random"
- Correctly leveraging ensemble methods which utilize multiple statistical and computational pattern recognition methods
 - PITF forecasting tournament; Bayesian model averaging
 - There are known and irreducible random elements in political behavior
- Upshot: you can't simply specify a desired rate of accuracy and assume by throwing sufficient money at the problem you will get there.
 - PITF and many other conflict forecasting models all converge to about 80% accuracy

Irreducible sources of error

- Specification error: no model of a complex, open system can contain all of the relevant variables;
- Measurement error: with very few exceptions, variables will contain some measurement error
 - presupposing there is even agreement on what the "correct" measurement is in an ideal setting;
 - Predictive accuracy is limited by the square root of measurement error: if your reliability is 80%, your accuracy can't be more than 90%
- Free will
 - Rule-of-thumb from our rat-running colleagues: "A genetically standardized experimental animal, subjected to carefully controlled stimuli in a laboratory setting, will do whatever it wants."
- Quasi-random structural error: Complex and chaotic deterministic systems behave as if they were random under at least some parameter combinations

Statistical challenges

- Rare events
 - Incorporate much longer historical time lines?—Schelling used Caesar's *Gallic Wars* to analyze nuclear deterrence
 - Calibration can be very tricky
- Analysis of event sequences, which are not a standard data type
- Causality
 - Oxford Handbook of Causation is 800 pages long
- Integration of qualitative and SME information
 - Bayesian approaches are promising but to date they have not really been used

Contemporary Technical Political Forecasting

- State Failures Project 1994
- Joint Warfare Analysis Center 1997
- FEWER [Davies and Gurr 1998]
- Various UN and EU forecasting projects
- Center for Army Analysis 2002-2005
- Swiss Peace Foundation FAST 2000-present
- Political Instability Task Force 2002-present
- DARPA ICEWS 2007-present
- Peace Research Center Oslo (PRIO) and Uppsala University UCDP political forecasting models

But these models don't work! Wired magazine tells me so!



PREVIOUS POST

NEXT POST

172

Tweet 🞗 +1 📊 Share

518

Pentagon's Prediction Software Didn't Spot Egypt Unrest

By Noah Shachtman 🖂 February 11, 2011 | 7:00 am | Categories: DarpaWatch

🖞 Like 🛛 🤤 Send 📑 497 people like this. Sign Up to see what your friends like.



In the last three years, America's military and intelligence agencies have spent more than \$125 million on computer models that are supposed to forecast political unrest. It's the latest episode in Washington's four-decade dalliance with future-spotting programs. But if any of these algorithms saw the upheaval in Egypt coming, the spooks and the generals are keeping the predictions very quiet.

The Forecaster's Trilogy

- Nassem Nicholas Taleb. *The Black Swan*[most entertaining]
- Daniel Kahneman. *Thinking Fast and Slow* [30 years of research which won Nobel Prize]
- Philip Tetlock. *Expert Political Judgment* [most directly relevant]

The Forecasting Zoo





Why should we care about ducks?



Size



Quantity



Variety



Suspicious behaviors

DARPA-World

كم as as as as as as as -de de han han han han han

NNT-World


The Forecasting Zoo







Owls: the narrative fallacy



- Tetlock: "Experts" predict only slightly better than chance; the more famous the expert, the lower the accuracy
- Hegel: the owl of Minerva flies only at dusk
- Taleb: seeking out narratives is an almost unavoidable cognitive function and it generates a dopamine hit
- "Loyalty!? Owls are only loyal to their stomachs!" (Kansas wildlife biologist commenting on owls as pets, per *Harry Potter* stories)

This is your brain on narratives





The Forecasting Zoo











Tigers

- You know they are out there
- You know they can hurt you
- But you don't know where they are

IC-World?



Questions?

Philip A. Schrodt Political Science Pennsylvania State University State College, PA 16802 Phone: 814-863-8978 Email: schrodt@psu.edu Project Web Site: http://eventdata.psu.edu

Open sources



Mainstream media



Internet and new social media

New social media

- The good
 - Widely available to elites
 - More or less uncensored
 - Should provide early information on changing sentiment prior to observing actual collective action
- The bad
 - No filters and mostly politically irrelevant: "Wanna getta pizza? ;)"
 - Easily manipulated by anyone—business, government, NGOs—who wants to go to the trouble of doing so
- The ugly
 - No standardization of content

New Social Media

Police tell high school students to disguise their identity on FaceBook

+

Students choose the first country they see in an alphabetical list.... Afghanistan

OMG! Jihadis are checking out the junior prom!

Open sources



Mainstream med



Internet and new social media

NSM: A Conjecture

 NSM actually amplify noise due to the premium on entertainment

Challenges in integrating models into decision-making

- Forecasting is hard (Tetlock)
- Probabilistic reasoning is hard (Kahneman)
- Statistics is new compared to deterministic modeling and is still changing, even at very fundamental levels
 - Frequentist vs Bayesian approaches
 - New approaches made possible by computational advances

Earliest date of a textbook that could be used in a contemporary class

- Geometry: 300 BCE
- Algebra: 1750
- Calculus and differential equations: 1820
- Statistics: ???
 - Current social science statistics books are mostly frequentist (hypothesis testing using significance levels) whereas most statistics departments are now Bayesian (estimating probability distributions of coefficients)
 - The topics in the introductory curriculum vary little from precomputer times

Event Data Generation Process



Human Geography: Open Source Data, Methods and Forecasting

Philip Schrodt Political Science and International Center for the Study of Terrorism Pennsylvania State University schrodt@psu.edu http://eventdata.psu.edu

4 September 2012

What's the buzz?

Why 2012 is not 2002



More general resource on this:

Manyika, James; Michael Chui, Jaques Bughin, Brad Brown, Richard Dobbs, Charles Roxburgh, Angela Hung Byers (May 2011). *Big Data: The next frontier for innovation, competition, and productivity.* McKinsey Global Institute.



Most of these can be located via Google, but contact me for additional info.

Also see the following sites:

http://www.paulhensel.org/data.html

http://www.isadiscussion.com/view/0/datasets.html

Political Data Collection: Penn State

- Correlates of War Project (Bennett, Palmer)
- Legislative Text (Monroe)
- Event Data (Schrodt)
- World Religions (Finke, Henderson, Maoz)
- Association of Religion Data Archives (Finke)
- International Center for the Study of Terrorism (Horgan, Bloom)
- NSF Integrated Graduate Research and Training: Big Data in the Social Sciences (Monroe)



The pictured cluster computer is actually at the NASA Center for Climate Simulation (http://www.nasa.gov/topics/earth/features/climate-sim-center.html); it looks cooler than ours.

Visualization center is one of four at Penn State: this one is the Immersive Construction (ICON) Lab, a partnership facility in the department of Architectural Engineering and the Computer Integration Construction research program: see

http://rcc.its.psu.edu/resources/facilities_partnerships/

Computing Power

Control Data Corporation 3600 (ca. 1965) 32 K (48-bit) RAM memory 1 processor

~1-million operations per second Output: line printer









Penn State High Performance Computing Facility 15 cluster computers 100 to 2000 2.66 Ghz processors in each cluster ~50 Gb RAM accessible to each processor 130 Tb disk space

4 interactive visualization rooms

Motorola Razr 16 Gb RAM memory Dual processor ~500-million operations per sec 540 x 860 color display



- So what? A research operation that ten years ago might require \$50,000 per researcher in software licenses can now assemble all of the required tools for free. Resources such as *stackoverflow* and *sourceforge* provide free crowd-sourced advice on how to use these, as well as sample code, and this information has been vetted by a very sophisticated and international professional community.
- A military analogy might be the transition from bronze to iron weapons. Bronze is relatively easy to work with---how do you say "user friendly" in Hittite?--- but the materials are scarce and expensive. Iron is more difficult to master but once you've done so, it is cheap, abundant, and produces a much better product.
- For the independent research community, this has been a complete game changer.







- Open source
- Widely used in all approaches to statistical analysis *and* pattern recognition
- CRAN library provides almost immediate access to new methods
 - events (event data aggregation)
 - tm (text miner)
- Robust scripting capabilities; easily interfaces with C/C++ when needed
- Skill set is widely available



This is from a proposal we currently have under consideration at the National Science Foundation Methodology, Measurement and Statistics program. It would replace several human coding projects with near-real-time automated systems, augmented by decentralized human input.



http://eventdata.psu.edu. Graphics generated using a R script.

- But wait! These data are coded with TABARI, and BBN, a gadzillion-dollar corporation, has definitively and irrevocably demonstrated that TABARI is crap, haven't they? And BBN has neither the incentive nor inclination to misrepresent this issue, right? Right??
- Well, perhaps, but before making up your mind, at least skim this memo: http://eventdata.psu.edu/papers.dir/Schrodt.responsetoBBN.pdf



Currently undergoing tests and quality checks, but we expect to release the data sometime in the fall of 2012. This was an informal collaboration between the University of Illinois and our NSF-funded coding operation at Penn State.



Data which appears to be missing along the equator was due to a bug in the visualization software.



http://www.acleddata.com/



http://www.acleddata.com/



http://www.acleddata.com/


http://emm.newsbrief.eu/overview.html



http://www.pcr.uu.se/research/UCDP/



Further reflections on the problems encountered in government contracting in these fields can be found at http://eventdata.psu.edu/7DS/7DS.Practitioners.chpt.pdf

Taleb, by the way, views the military as one of the few professional groups with a realistic approach to dealing with uncertainty, albeit his baseline is Wall Street traders and academic economists, and he is dealing with the elite uniformed military he encountered at the likes of the Aspen Institute, not the larger defense consulting community.



- Reductionist models: expected utility, game theory, systems dynamics, and agent-based models.
- See Kahneman and Taleb for about 800 pages on why the simplifying assumptions of expected utility and game theory doom those approaches to irrelevance as models of human behavior. Even in methodologicallyconservative academic economics, these approaches have been largely abandoned in favor of newer evolutionary approaches and evidence-based "behavioral economics" models.
- Systems dynamics were doomed by the [re-]discovery of chaos theory in the 1970s (as well as the conspicuous failure of a large number of very expensive models): Add a quadratic term, and depending on coefficient choices your model is potentially chaotic: this effect is so simple it can be demonstrated in Excel. Keep everything linear and you've just got a model that any undergraduate who has gotten halfway through a differential equations course can solve. Value added: zero. Outside the defense research community, systems dynamics was completely abandoned by the early 1980s.

The jury is still out on agent-based models: The problem is calibration and metrics for assessing accuracy, though there may be solutions to this.



Proprietary "Give me \$10-million to write you some software which you can't examine and I plan to license back to you for \$10,000/seat/year" models as frisbee-catching dogs? Yeah, sounds about right.

Apologies to those who have trained frisbee-catching dogs.



Frequentist models useless for prediction:

Michael D. Ward, Brian D. Greenhill, and Kristin M. Bakke. The perils of policy by p-value: Predicting civil conflicts. Journal of Peace Research, 47(5), 2010.

Additional criticisms of frequentist practice in the social sciences: http://eventdata.psu.edu/7DS/Schrodt.7Sins.APSA10.pdf

Statistical challenges

- Systematically dealing with measurement error and missing values rather than assuming "missing at random"
- Correctly leveraging ensemble methods which utilize multiple statistical and computational pattern recognition methods
 - · PITF forecasting tournament; Bayesian model averaging
 - There are known and irreducible random elements in political behavior
- Upshot: you can't simply specify a desired rate of accuracy and assume by throwing sufficient money at the problem you will get there.
 - PITF and many other conflict forecasting models all converge to about 80% accuracy





Thomas Schelling. Arms and Influence. 1966

Contemporary Technical Political Forecasting

- State Failures Project 1994
- Joint Warfare Analysis Center 1997
- FEWER [Davies and Gurr 1998]
- Various UN and EU forecasting projects
- Center for Army Analysis 2002-2005
- Swiss Peace Foundation FAST 2000-present
- Political Instability Task Force 2002-present
- DARPA ICEWS 2007-present
- Peace Research Center Oslo (PRIO) and Uppsala University UCDP political forecasting models

Davies, John L. and Ted R. Gurr, eds. 1998. Preventive Measures: Building Risk Assessment and Crisis Farly Warping, Lanham, MD: Rowman and Littlefield

Crisis Early Warning. Lanham, MD: Rowman and Littlefield.

Esty, Daniel C., Jack A. Goldstone, Ted R. Gurr, Pamela Surko, and Alan N. Unger. 1995. State

Failure Task Force Report. McLean, VA: Science Applications International Corporation. Esty, Daniel C., Jack A. Goldstone, Ted R. Gurr, Barbara Harff, Marc Levy, Geoffrey D.

Dabelko, Pamela Surko, and Alan N. Unger. 1998. State Failure Task Force Report: Phase II

Findings. McLean, VA: Science Applications International Corporation.



- Source: http://www.wired.com/dangerroom/2011/02/pentagon-predict-egyptunrest/
- Further into the story, Shachtman acknowledges that ICEWS—still in the developmental stages in 2011 and focusing exclusively on Asia—wasn't even *monitoring* the Middle East. Might have had something to do with ICEWS not predicting the Arab Spring, ya think?
- Shachtman, for whatever reason, seems to have a beef about technical forecasting models in general and ICEWS in particular. On the positive side, at least from the perspective of publicity-shy ICEWS, he has about 50% of his facts wrong about the system.

See also: http://www.wired.com/dangerroom/2007/11/lockheed-peers/

The Forecaster's Trilogy

- Nassem Nicholas Taleb. *The Black Swan* [most entertaining]
- Daniel Kahneman. *Thinking Fast and Slow* [30 years of research which won Nobel Prize]
- Philip Tetlock. *Expert Political Judgment* [most directly relevant]



Black swans: very low probability, high consequence events (Taleb) Rubber ducks: conventional events (Schrodt)





- Problem: DARPA is utterly clueless about social science research and views it as an exotic technology which should only be used in black swan situations. Meanwhile the rest of the world speeds ahead of them.
- Norway with 1/60th the US population spends 6-times as much on quantitative international conflict research. Why?—perhaps they can't afford not to.
- The current (August 2012) Minerva BAA explicitly requests funding for evermore game theory models. Smooth-bore muskets? Good enough for Napoleon; good enough for me.



NNT = Nassem Nicholas Taleb, as he refers to himself. This 15/85 split is based on Taleb's experience as a Wall Street trader, where he succeeded to a point where he could give \$100 tips to cabbies, as well as writing books. While based on experience in finance, he notes that this is not a bad balance generally. 15% attention to black swans is generally sufficient to keep one out of "sucker" territory.





Reflections on the narrative fallacy:

- Joll, James. The Origins of the First World War (1984) discusses [at least] 14 theories as to why WWI was "inevitable." http://en.wikipedia.org/wiki/Causes_of_World_War_I lists 46 books on the topic. However, as late as July 1914 there was little concern expressed in the European press about a major war, and attempts to duplicate the outbreak in simulations in the 1960s almost always failed.
- Then consider the summer 2012 USA Network cable television mini-series Political Animals. It focuses on the travails of a female Secretary of State. Okay, that's plausible enough. But adding to the drama, she's an experienced political actor in her own right, and almost defeated the President who appointed her in the primaries! But it gets wilder. She's actually the wife of a *former* President! And---this is Hollywood---that former President is a notorious womanizer!
- Give me a break! Next thing know you, they will be making the President an African-American with an Arabic name who was raised in Hawaii by his Kansas-born grandmother and ends up going to Harvard Law School and teaching at the University of Chicago after working on the Chicago's South Side as a community organizer.
- But they wouldn't dare push reality that far: in the mini-series, the President is a white guy.



I'm going to see if I can raise money on KickStarter to place a bunch of these ads in the Pentagon and Foggy Bottom metro stops...



Tigers



- You know they are out there
- You know they can hurt you
- But you don't know where they are





When I just saw the smaller version of that metal-band image, I thought it said "Children of Boredom." Yeah, that's social media alright...



New Social Media Police tell high school students to disguise their identity on FaceBook + Students choose the first country they see in an alphabetical list... Afghanistan OMG! Jihadis are checking out the junior prom!

Source: http://www.danah.org/papers/TakenOutOfContext.pdf

"danah boyd" [sic] is now a senior researcher at Microsoft: http://en.wikipedia.org/wiki/Danah_boyd



NSM is to the IC in the 2010s as crack cocaine was to urban ghettoes in the 1980s

NSM: A Conjecture

NSM actually amplify noise due to the premium on entertainment

Challenges in integrating models into decision-making

- Forecasting is hard (Tetlock)
- Probabilistic reasoning is hard (Kahneman)
- Statistics is new compared to deterministic modeling and is still changing, even at very fundamental levels
 - Frequentist vs Bayesian approaches
 - New approaches made possible by computational advances

Earliest date of a textbook that could be used in a contemporary class

- Geometry: 300 BCE
- Algebra: 1750
- Calculus and differential equations: 1820
- Statistics: ???
 - Current social science statistics books are mostly frequentist (hypothesis testing using significance levels) whereas most statistics departments are now Bayesian (estimating probability distributions of coefficients)
 - The topics in the introductory curriculum vary little from precomputer times