# Automated Production of High-Volume, Near-Real-Time Political Event Data

Philip A. Schrodt

Pennsylvania State University

schrodt@psu.edu

# What do these events have in common?

- 1941 Pearl Harbor attack
- 1950 Chinese intervention in Korea
- 1962 Cuban missile crisis
- 1973 Middle East War
- 1979 Iranian Revolution
- 1989 End of Soviet-supported regimes in Eastern Europe
- 1994 Rwanda genocide
- 1998 Indian nuclear tests
- 2001 Al-Qaeda attack on U.S.

# Surprise!

# Surprise!

- Strategic surprise: unexpected events with major consequences
- This despite the fact that people were supposed to be watching

# Surprise!

- Strategic surprise: unexpected events with major consequences
- This despite the fact that people were supposed to be watching
  - $50-billion to $80-billion of watching

Protecting yourself from political surprise…

Some approaches…

Protecting yourself from political surprise…

Some approaches…

"The further you look back, the further you can see ahead"—Winston Churchill

# Early technical forecasting models

- Divination model of sheep liver

- Source: Babylonia, ca. 600 BC

# Early technical forecasting models
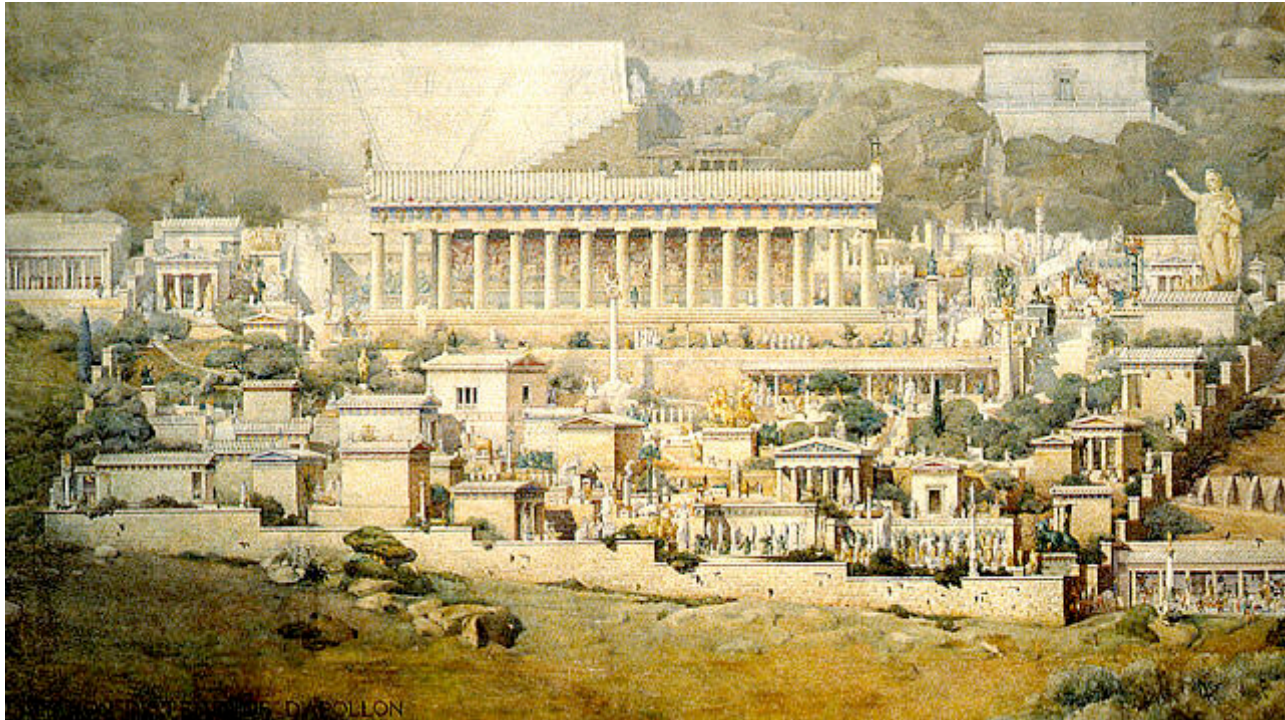
- Divination model of sheep liver

- Source: Babylonia,  ca. 600 BC
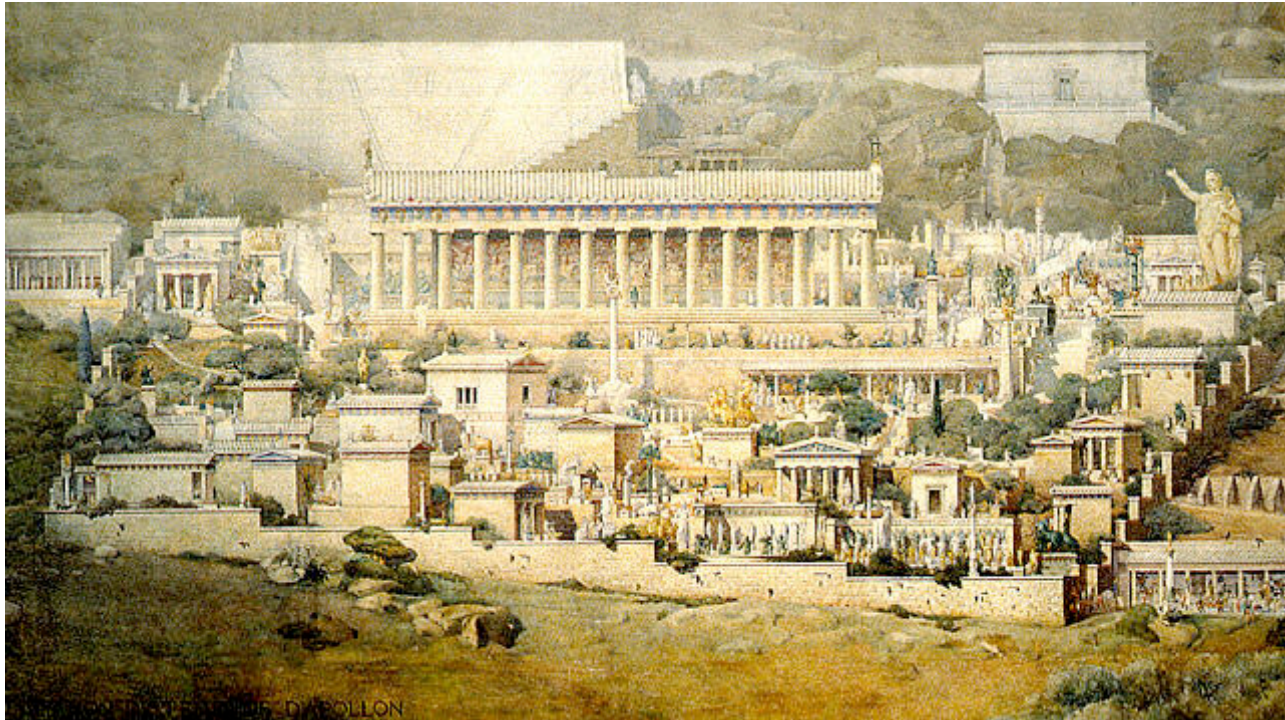
- Persian conquest of Babylonia
  539 BCE

# Temple of Apollo at Delphi

# Temple of Apollo at Delphi



Methodology:

- Inhale hallucinogenic vapors
- Make ambiguous predictions: "A mighty kingdom will fall"

Modern solution: the ubiquitous SME

"Subject matter expert"

# SMEs: Self-Perception

# SMEs: Reality

# SMEs: Real Reality



The New York Times

**Politics**

WORLD | U.S. | N.Y. / REGION | BUSINESS | TECHNOLOGY | SCIENCE | HEALTH | SPORTS | OPINION

POLITICS HOME | INSIDE CONGRESS | HOUSE | SENATE | GOVERNORS

## Hiding Details of Dubious Deal, U.S. Invokes National Security

By ERIC LICHTBLAU and JAMES RISEN
Published: February 19, 2011

WASHINGTON — For eight years, government officials turned to Dennis Montgomery, a California computer programmer, for eye-popping technology that he said could catch terrorists. Now, federal officials want nothing to do with him and are going to extraordinary lengths to ensure that his dealings with Washington stay secret.

The Justice Department, which in the last few months has gotten protective orders from two federal judges keeping details of the technology out of court, says it is guarding state secrets that would threaten national security if disclosed. But others involved in the case say that what the government is trying to avoid is public embarrassment over evidence that Mr. Montgomery bamboozled federal officials.

Dennis Montgomery

A onetime biomedical technician with a penchant for gambling, Mr. Montgomery is at the center of a tale that features terrorism scares, secret White House briefings, backing from prominent Republicans, backdoor deal-making and fantastic-sounding computer technology.

**Multimedia**

Interactive Feature
Boom Time to Bankruptcy

Interviews with more than two dozen current and former officials and business associates and a review of documents show that Mr. Montgomery and his associates received more than $20 million in government contracts by claiming that software he had developed could help stop Al Qaeda's

# "McChrystal's Hairball"



Afghanistan Stability / COIN Dynamics

WORKING DRAFT – V3

# Problems with SMEs

- Like diamonds, expensive and rare
- Affected by fatigue and boredom, cognitive dissonance and consulting fees
  - "You want it real bad, you're gonna get it real bad…"
  - Usually (not always) considerably more entertaining and articulate than computer models
- SMEs with native language and cultural expertise may have axes to grind…
  - See N. Machiavelli, 1513 on trusting exiles
- Accuracy is only slightly better than chance on hard problems
  - The more visible the SME, the less accurate the record

# Technical forecasting as the alternative

Mix in proper proportions

- Political science theory

- Data

- Statistical Models

- Computing power

Result: models accurate in the 75% - 85% range

# Cantor-Fitzgerald forecasting algorithms

Cantor-Fitzgerald forecasting algorithms

+

All the sports data Cantor-Fitzgerald can afford

(all the sports data that exists)

Cantor-Fitzgerald forecasting algorithms

+

All the sports data Cantor-Fitzgerald can afford

(all the sports data that exists)

+

All the computing power Cantor-Fitzgerald can afford

(a lot of computing power)

Cantor-Fitzgerald forecasting algorithms

+

All the sports data Cantor-Fitzgerald can afford

(all the sports data that exists)

+

All the computing power Cantor-Fitzgerald can afford

(a lot of computing power)

A sports betting system

# Bling!

# Cantor-Fitzgerald MIDAS System

# Integrated Conflict Early Warning System

- Funded by DARPA Information Processing Techniques Office [unclassified]
- Funding at $40-million for 2007-2011
  - Largest quantitative conflict analysis project since the 1970s
- Objective is real-time forecasting of indicators of political instability in Asia with 6-24 month leads, 70%-80% accuracy
- Machine-coded event data has developed as the core methodology
  - Data covers 1997-present with 8.5-million stories from 27 sources

# Political Instability Task Force

- Multi-agency, including State, Census, DoD, but primarily CIA [unclassified]
- Initially established in 1995 as "State Failures Project" in response to former Yugoslavia, Somalia and Rwanda
- Global forecasting of state failure, democratic-authoritarian transitions, mass killings and other indicators at 6-24 month leads, models have demonstrated 70%-80% accuracy
- Funds a number of major data collection efforts
- Other participating academic institutions include Harvard, Stanford, Columbia, Michigan, George Mason

# Political Instability Task Force (AJPS 2010)

**TABLE 2   Out-of-Sample Prediction Exercise for Observed Onsets of Instability, 1995–2004**

### A. Countries That Had Instability Onsets, 1995–2004. Quintile/decile in model score rankings based on 2-yr. prior data

| Year | Top Decile | Second Decile | Second Quintile | Third Quintile |
|---|---|---|---|---|
| 1995 | Armenia, Comoros | Belarus | | |
| 1996 | Albania, Niger, Zambia | | Nepal | |
| 1997 | Cambodia, Congo-Brazz. | | | |
| 1998 | Guinea-Bissau, Lesotho | | | Serbia/Montenegro |
| 1999 | Ethiopia, Haiti | | | |
| 2000 | | Solomon Ils., Guinea* | | |
| 2002 | Cote d'Ivoire | | | |
| 2003 | Central African Republic | | | |
| 2004 | Iran* | Yemen* | | Thailand* |

### B. Tabulation of All Country-years, 1995–2004. Model estimates based on censored data, using only sample data from prior to year of forecast (countries w/population over 500,000, no ongoing conflict, at least two years old)

| | Countries with Instability in $t + 2$ | Countries Remaining Stable |
|---|---|---|
| Predicted for Instability (Top Quintile) | 18 | 233 |
| Predicted for Stability (Not Top Quintile) | 3 | 992 |
| N = 1,246 Percent Classed Correctly | 85.7% | 81.0% |

Number of instability onsets, 1995–2004: 21. Number of instability onsets in top quintile of model scores: 18 (86%).
*Cases added to the problem set in 2005 update.

# Example: 18 December 2007

BAGHDAD — Iraqi leaders criticized Turkey on Monday for bombing Kurdish militants in northern Iraq with airstrikes that they said had left at least one woman dead.

The Turkish attacks in Dohuk Province on Sunday — involving dozens of warplanes and artillery — were the largest known cross-border attack since 2003. They occurred with at least tacit approval from American officials. The Iraqi government, however, said it had not been consulted or informed about the attacks.

Massoud Barzani, leader of the autonomous Kurdish region in the north, condemned the assaults as a violation of Iraqi sovereignty that had undermined months of diplomacy. "These attacks hinder the political efforts exerted to find a peaceful solution based on mutual respect."

*New York Times*, 18 December 2007

http://www.nytimes.com/2007/12/18/world/middleeast/18iraq.html?_r=1&ref=world&oref=slogin (Accessed 18 December 2007)

# TABARI Coding

BAGHDAD — Iraqi leaders criticized Turkey on Monday for bombing Kurdish militants in northern Iraq with airstrikes that they said had left at least one woman dead.

# TABARI Coding: Verb

BAGHDAD — Iraqi leaders criticized Turkey on Monday for bombing Kurdish militants in northern Iraq with airstrikes that they said had left at least one woman dead.

Event Code: 111

# TABARI Coding: Actors

BAGHDAD — Iraqi leaders criticized Turkey on Monday for bombing Kurdish militants in northern Iraq with airstrikes that they said had left at least one woman dead.

Event Code: 111

Source: IRQ

Target: TUR

# TABARI Coding: Agents

BAGHDAD — Iraqi leaders criticized Turkey on Monday for bombing Kurdish militants in northern Iraq with airstrikes that they said had left at least one woman dead.

Event Code: 111

Source: IRQ  GOV

Target: TUR

# TABARI Coding: Verb [2]

BAGHDAD — Iraqi leaders criticized Turkey on Monday for bombing Kurdish militants in northern Iraq with airstrikes that they said had left at least one woman dead.

Event Code: 111

Source: IRQ   GOV

Target: TUR

Event Code: 223

# TABARI Coding: Actors [2]

BAGHDAD — Iraqi leaders criticized Turkey on Monday for bombing Kurdish militants in northern Iraq with airstrikes that they said had left at least one woman dead.

Event Code: 111

Source: IRQ   GOV

Target: TUR


Event Code: 223

Source: TUR

Target: IRQKRD

# TABARI Coding: Agents [2]

BAGHDAD — Iraqi leaders criticized Turkey on Monday for bombing Kurdish militants in northern Iraq with airstrikes that they said had left at least one woman dead.

Event Code: 111

Source: IRQ   GOV

Target: TUR

Event Code: 223

Source: TUR

Target: IRQKRD    REB

# Event Model: Core Innovation

- Once calibrated, real-time event forecasting models can be run **entirely** without human intervention
  - Web-based news feeds provide a rich multi-source flow of political information in real time
  - Statistical models can be run and tested automatically, and are 100% transparent
- In other words, for the first time in human history—quite literally—we have a system that can provide real-time measures of political activity without any human intermediaries

# Event Data in the [original] DARPA period

- Nation-state orientation; most analysis dealt with Cold War major power relations

- Global coverage

- Human coding

- Source texts were major Western newspapers

- Statistical models were relatively simple

# Contemporary Event Data

- Substate and nonstate actors; most analysis deals with protracted conflicts

- Single-conflict and regional coverage

- Automated coding

- Source texts are from wire services (Reuters, AFP)

- Statistical models are very complex

# ICEWS Event Data

- 30-gigabytes of text from Lexis-Nexis
- 25 sources
- 8-million stories
- 26-million sentences
  - Only first four sentences coded in each story
- 3-million events
- Generally two orders of magnitude greater than any prior event coding effort

# Event Data Generation Process

# Raven Phase 2 Functional View

# Event Data Generation Process

# ICEWS "Events of Interest"

Domestic Political Crisis—Significant opposition to the government, but not to the level of rebellion or insurgency (for example, power struggle between two political factions involving disruptive strikes or violent clashes between supporters)

Rebellion—Organized opposition where the objective is to seek autonomy or independence

Insurgency—Organized opposition where the objective is to overthrow the central government

Ethnic/Religious Violence—Violence between ethnic or religious groups that is not specifically directed against the government

International Crisis—Conflict between two or more states or elevated tensions between two or more states that could lead to conflict

# Phase 1 Results: LM-ATL Out-of-Sample Results (DARPA Chart)



- Exceeds metrics for the maximum intensity index and 3 instability events: Rebellion, Insurgency, and Ethnic/Religious Violence – Passes Phase 1 gates
- By integrating improved versions of best of breed models from multiple perspectives, team achieves more accurate, precise forecasts than any one model alone

# Components of a forecast

Pie chart of impactful variables for country, EOI, time

- Impact power
- Model weight
- correlation

Variable 1
Variable 2
Variable 3
Variable 4
Variable 5
Other

**Relative Normalized Ranges**

Normalized Range

All Country Max
Individual Country Max
Individual Country Min
All Country Min

Variable 1
Variable 2
Variable 3
Variable 4
Variable 5

**Variable Values Impact on Prob(EOI)**

Impact increases P(EOI)

Impact decreases P(EOI)

Impact on P(EOI)

- Potential Increase to P(EOI)
- Potential Decrease to P(EOI)

Variable 1
Variable 2
Variable 3
Variable 4
Variable 5

**Prob(EOI) - Variable 1 Effect**

Current value

P(EOI)

Variable 1

# Challenges

- Lexis-Nexis
  - It was not designed for this sort of thing
  - Money is not the issue; the institution is the issue
- Duplicate detection
  - Considerably complicated by multiple sources
- Irrelevant story detection: at least a third of the downloads
- Dictionary development
  - 15,000 verb phrases are mostly adequate
  - Actor dictionaries—currently around 15,000 entries—were a major challenge
  - CountryCodes.xml generic dictionary shared with NSF/MID work
- TABARI in parallel coding mode works at 70,000 events per second
  - Parallelism on the cheap—just split the files

# News & Other Data Sources (Phase 1)

- TABARI open source event data coding tool
  - 6.7M news stories from 75+ sources
  - 253M lines of text
  - 30 dictionaries, 20K entries
  - CAMEO action taxonomy
  - complementing with AeroText in Phase 2
- Country/State data
  - 16+ sources
- SME interviews for agent-based country models

*It is estimated that this is the largest automated event coding project to date. Enabled by end-to-end automated process.*



| As of Jun 2008 | | Australia | Bangladesh | Bhutan | Burma (Myanmar) | Cambodia | China | Comoros | Fiji | India | Indonesia | Japan | Korea, North | Korea, South | Laos | Madagascar | Malaysia | Mauritius | Mongolia | Nepal | New Zealand | Papua New Guinea | Philippines | Russia | Singapore | Solomon Islands | Sri Lanka | Taiwan | Thailand | Vietnam |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Int'l | BBC World Monitoring Service* | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X |
| | The New York Times | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X |
| | The Associated Press | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X |
| | Japan Economic News Wire | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X |
| | United Press International | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X |
| | Inter Press Service | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X |
| | Agence France Presse | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X |
| | Xinhua General News Service | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X |
| | Deutsche Presse Agentur | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X | X |
| Local Sources | Bangkok Post | | | | X | X | | | | X | | | | X | | | | | | | | | | | | X | | | X | X |
| | Central News Agency- Taiwan | | | | X | X | | | | | X | X | | | | | | X | | | | | | | | | X | | |
| | India Today | | X | X | | X | | | | X | | | | | | | | X | | | | | | | | X | | |
| | The Edge Malaysia | | | | X | | | | | X | | | | X | | | | X | | | | | | X | | | X | X |
| | The Edge Singapore | X | | | X | X | | | | X | | | | X | | | | X | | X | X | | X | | | | X | X |
| | The Jakarta Post | X | | | X | X | | | | X | | | | X | | | | X | | X | X | | X | | | | X | X |
| | The Nation (Thailand) | X | | | X | X | | | | X | | | | X | | | | X | | X | | X | | | | X | X |
| | The Nation (Pakistan) | | X | X | | X | | | | X | | | | | | | X | | | | | | | X | | |
| | The Pakistan Newswire | | X | X | | X | | | | X | | | | | | | X | | | | | | | X | | |
| | The Statesman- India | | X | X | | X | | | | X | | | | | | | X | | | | | | | X | | |
| SMEs | UPENN Questionaire | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | IDI Leaders Questionaire | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |

Being Purchased or Developed
Acquired
X  Encoded (multiple times)

*Operate around the clock to monitor more than 3,000 radio, TV, press, internet and news agency sources, translating from up to 100 languages.

# [Near-] Real Time Coding

- ICEWS is currently doing monthly updating
- Experimental system at Kansas did daily coding based on web site news feeds from Reuters and UPI
  - This is be implemented and enhanced under NSF funding at UT/Dallas (Brandt)
  - This is will linked to a real-time prediction module

# Expansion to Global Coverage

- Projecting an additional 25M stories, 2000-2011
  - International sources appear to provide most of the coverage in English
- Intake of a existing political entity lists such as rulers.org, CIA World Factbook, Joshua Project
- Rapid development of remaining actor dictionaries using NER tools
- XML-based Hierarchical dictionaries: Move entity information to a database and generate specific codes as needed
- Objective: complete project by August 2011

# Third Generation Coder using Open Source NLP Tools

- Full parsing using the TreeBank parse tree format
- Named-entity recognition/resolution software
- Parts of speech markup for noun-verb disambiguation
- Stemming: Porter stemmer
- WordNet synsets
- Sentence delineation
- GeoNet Names Server for geographical place names
- Regular expressions for dictionary syntax

# JABARI-NLP Performance

- TABARI baseline: 56% precision, 54% recall

- Add Open-NLP Penn TreeBank parser:
  68% precision, 35.4% recall

- Add GATE-Annie noun phrase synonyms,
  pronoun coreferencing, and default location
  agent resolution:
  77% precision, 66.5% recall

# Additional needed tasks

- Duplicate detection and content classification
  - Google News and European Media Monitor
- Hierarchical dictionaries
  - Move entity information to a database and generate specific codes as needed
  - Automated dictionary updating based on dense sources such as rulers.org
- Move entire system to the web
  - Separate graphical and processor-intensive components
- Sophisticated error detection/correction

# Unsolicited career advice

# Unsolicited career advice

- Learn computer programming
- Eric Drexler: "data-driven science"
- *Fourth Paradigm: Data Intensive Scientific Discovery* (Hey, Tansley and Toole, eds.)

# Statistics Packages :: Cars

# Statistics Packages :: Cars

Stata

# Statistics Packages :: Cars

Stata



SAS

# Statistics Packages :: Cars

Stata

R

SAS

# Statistics Packages :: Cars

Stata

R

SAS

SPSS

# NPL Components :: Smart Phones

# NLP… it's all components

# Evolution of Statistical Training in Political Science

# Evolution of Statistical Training in Political Science

- 1985: Objective is to get grad students to take advanced courses in econometrics
  - Hint: deadends— SEM, co-integration models

# Evolution of Statistical Training in Political Science

- 1985: Objective is to get grad students to take advanced courses in econometrics
  - Hint: deadends— SEM, co-integration models
- 2010: Objective is to develop advanced techniques that will be adopted in other disciplines
  - Imai, Sekhon, Fowler, Gill, King
  - Hint: Medical schools are rumored to pay better than the liberal arts

# The Transition We Need in Programming

- Old model:
  - "We'll just hire a programmer because that will be more efficient than doing it ourselves"

# The Transition We Need in Programming

- Old model:
  - "We'll just hire a programmer because that will be more efficient than doing it ourselves"
- Reality
  - Computer science departments and ExxonMobil can't find enough programmers either
  - You take a serious efficiency hit in trying to explain what you want done
  - You may take a serious efficiency hit in not doing the task in the best way—NLP (and statistics) are specialized subfields
  - *Jurassic Park* vs. computational linguistics

# We need…

- Some formal training in core algorithms and data structures
  - Not just AP Java
- Rapid development scripting languages: perl and Python
- Lingua franca (and GUI): Java
- High performance: C/C++

# Three more things I'm not supposed to say…

# Three more things I'm not supposed to say...

- From a single source...

# Bacon!

# Bacon!



I delicious meats

Tactical Bacon
smoke flavor added
Fully Cooked!

NET WT. 9 OZ

Uncle Oinker's
savory
Bacon MINTS

Net Weight 0.7oz (20g)  100 Mints

No, not that kind of Bacon...

# Bacon!



Francis Bacon [1561 – 1626]

# Relevant scientific principles from Sir Francis

# Relevant scientific principles from Sir Francis

- The entire scientific method

# Relevant scientific principles from Sir Francis

- The entire scientific method
  - With additional help from Galileo and Descartes
  - Won't want to discard all that now, would we?

# Relevant scientific principles from Sir Francis

- Knowledge comes from induction, not just theory

# Relevant scientific principles from Sir Francis

- Knowledge comes from induction, not just theory
  - What's wrong with theory??? I'm not a real political scientist unless I have many, many, many theories…

# Yeah, right...



Theories of INTERNATIONAL POLITICS AND ZOMBIES

Daniel W. Drezner

# Or this?...

# Relevant scientific principles from Sir Francis

- Knowledge comes from induction, not just theory
  - "If you've got too much data, you need better theory' if you've got too much theory, you need better data"
  - When theory trumps data, you have Scholasticism

# Relevant scientific principles from Sir Francis

- Knowledge comes from induction, not just theory
  - "If you've got too much data, you need better theory' if you've got too much theory, you need better data"
  - When theory trumps data, you have Scholasticism
  - Which dominated the academic world for 250 years after Bacon and Descartes
    - ° Which is hardly reassuring

# Relevant scientific principles from Sir Francis

- A theory is only scientific if it is predictive
  – With help from the logical positivists, in particular Quine and Hempel

# Relevant scientific principles from Sir Francis

- A theory is only scientific if it is predictive
  - With help from the logical positivists, in particular Quine and Hempel
  - Predictive power is the only feature that distinguishes astronomy from astrology
    - Really…this is not just a cheap shot

# Relevant scientific principles from Sir Francis

- A theory is only scientific if it is predictive
  - With help from the logical positivists, in particular Quine and Hempel
  - Predictive power is the only feature that distinguishes astronomy from astrology
    - Really…this is not just a cheap shot
  - Focus on prediction and you will also wean yourself from frequentism
    - p-value frequentism and the logical-hypothetical method are mutually inconsistent

# Relevant scientific principles from Sir Francis
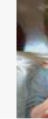
- Science should serve the public good

# Relevant scientific principles from Sir Francis

- Science should serve the public good
  - At least think about trying to developing statistical methods for applied forecasting

# Relevant scientific principles from Sir Francis

- Science should serve the public good
  - At least think about trying to developing statistical methods for applied forecasting
  - Maj Gen Michael Flynn, *Fixing Intel: A blueprint for making intelligence relevant in Afghanistan*
  - They need all the help they can get…

# Hiding Details of Dubious Deal, U.S. Invokes National Security

By ERIC LICHTBLAU and JAMES RISEN
Published: February 19, 2011

WASHINGTON — For eight years, government officials turned to Dennis Montgomery, a California computer programmer, for eye-popping technology that he said could catch terrorists. Now, federal officials want nothing to do with him and are going to extraordinary lengths to ensure that his dealings with Washington stay secret.

The Justice Department, which in the last few months has gotten protective orders from two federal judges keeping details of the technology out of court, says it is guarding state secrets that would threaten national security if disclosed. But others involved in the case say that what the government is trying to avoid is public embarrassment over evidence that Mr. Montgomery bamboozled federal officials.

Dennis Montgomery

A onetime biomedical technician with a penchant for gambling, Mr. Montgomery is at the center of a tale that features terrorism scares, secret White House briefings, backing from prominent Republicans, backdoor deal-making and fantastic-sounding computer technology.

**Multimedia**

Interactive Feature

Boom Time to Bankruptcy

Interviews with more than two dozen current and former officials and business associates and a review of documents show that Mr. Montgomery and his associates received more than $20 million in government contracts by claiming that software he had developed could help stop Al Qaeda's

# 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

# LM ATL Results by Model



Average distance, in quarters and over 29 countries, between probability prediction and ground truth vector for Rebellion, 2005-2006

Lower value = less error

LM ATL | UW | LC | UPenn | SAE L | UK | SAE B

**Aggregating model scores with a learned Bayesian network outperforms any one other model**

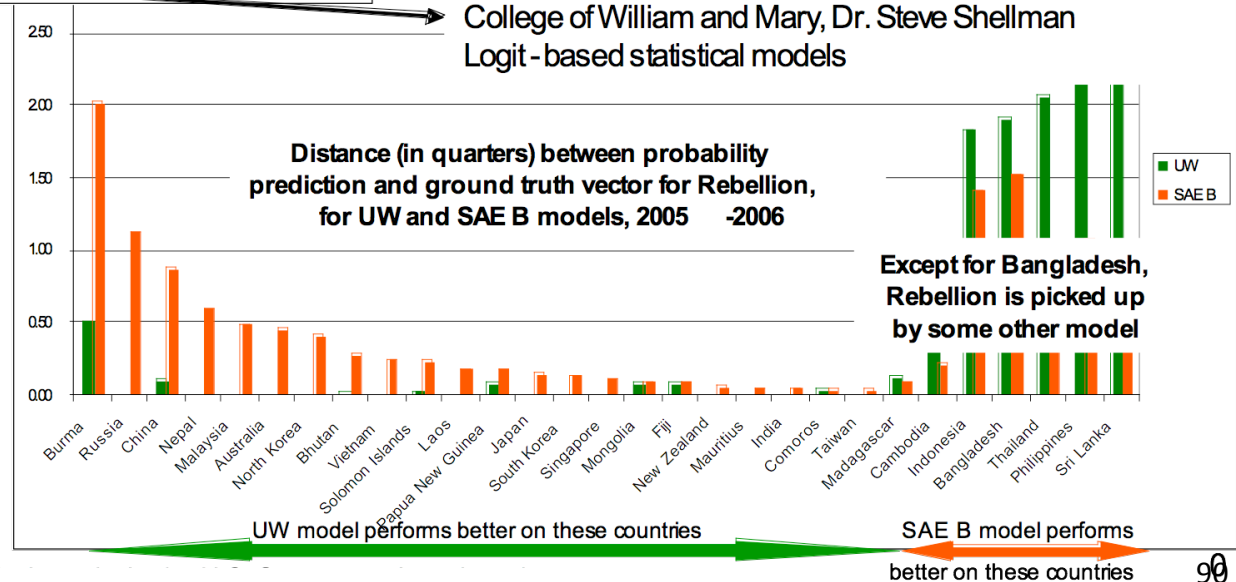– Because different models cover the EOIs and countries with varying levels of performance

College of William and Mary, Dr. Steve Shellman
Bayesian statistical model

University of Kansas, Dr. Phil Schrodt
College of William and Mary, Dr. Steve Shellman
Logit-based statistical models

Aggregation model

University of Penn, Dr. Ian Lustic and Dr. Barry Silverman
Agent-base models (only 6 and 4 "hard" countries resp.)

University of Washington, Dr. Michael Ward
Geo-spatial statistics applied to trade ties, flow of people, social similarity

Distance (in quarters) between probability prediction and ground truth vector for Rebellion, for UW and SAE B models, 2005-2006

Except for Bangladesh, Rebellion is picked up by some other model

UW
SAE B

Burma, Russia, China, Nepal, Malaysia, Australia, North Korea, Bhutan, Vietnam, Solomon Islands, Laos, Papua New Guinea, Japan, South Korea, Singapore, Mongolia, Fiji, New Zealand, Mauritius, India, Comoros, Taiwan, Madagascar, Cambodia, Indonesia, Bangladesh, Thailand, Philippines, Sri Lanka

UW model performs better on these countries

SAE B model performs better on these countries

90

# ICEWS Evaluation Criteria

- Accuracy $= \dfrac{\text{\# of correct predictions}}{\text{\# of predictions made}}$

- Recall $= \dfrac{\text{\# of correctly predicted conflicts}}{\text{\# of conflicts that occurred}}$

- Precision $= \dfrac{\text{\# of correctly predicted conflicts}}{\text{\# of conflicts predicted to occur}}$

# Contemporary Event Data

- Substate and nonstate actors; most analysis deals with protracted conflicts

- Single-conflict and regional coverage

- Automated coding

- Source texts are from wire services (Reuters, AFP)

- Statistical models are very complex