Data Mining At the Crossroads: Successes, Failures, and Learning from Them

> Haym Hirsh Department of Computer Science Rutgers University Division of Information and Intelligent Systems U.S. National Science Foundation

What is Data Mining?

For the purposes of my presentation:

Data Mining = The extraction of useful information from data

(I.e., Data Mining broadly construed)

• Web search

- Web search
- Spam filtering

- Web search
- Spam filtering
- Recommender systems

- Web search
- Spam filtering
- Recommender systems
- Machine translation

- Web search
- Spam filtering
- Recommender systems
- Machine translation
- Massive data clusters

- Web search
- Spam filtering
- Recommender systems
- Machine translation
- Massive data clusters
- Conferences like this one: Participation by people in diverse, previously disjoint subfields (databases, machine learning, statistics, etc.)

- Web search
- Spam filtering
- Recommender systems
- Machine translation
- Massive data clusters
- Conferences like this one: Participation by people in diverse, previously disjoint subfields (databases, machine learning, statistics, etc.)
- Benchmark datasets

Copyright © 2007 Haym Hirsh

• Socio-Political

🕤 🕤 👻 🗧 Kitp:	://query.nytimes.com/gst/fullpage.html?sec=technology&res=9E0DEEDC103BF933A25754C0A9629C8B6 ist for Purge of Voters Proves Flawed - New Y	53						
Pop-ups temporarily a	allowed. To always allow pop-ups from this site, click here							nanannesa
	HOME PAGE MY TIMES TODAY'S PAPER VIDEO MOST POPULAR TIMES TOPICS TimesSelect Free 14-Day Trial Welcome, telephanto Membe							
	Tuesday, August 14, 2007 Technolo	gy						Search
	WORLD U.S. N.Y. / REGION BUSINESS TECHNOLOGY SCIENCE HEALTH	SPORTS	OPINION	ARTS	STYLE	TRAVEL	JOBS	REAL
	CIRCUITS CAMCORDERS CAMERAS CELLPHONES COMPUTERS HANDHELDS HO		MUSIC PER		S WI-FI	DOWNLOAD	s	
	 By PORD PESSENDEN Profilement July 10, 2004 Florida election officials used a flawed method to come up with a listing of people believed to be convicted felons, a list that they are recommending be used to purge voter registration rolls, state officials acknowledged yesterday. As a result, voters identifying themselves as Hispanic are almost completely absent from that list. Of nearly 48,000 Florida residents on the felon list, only 61 are Hispanic. By contrast, more than 22,000 are African-American. About 8 percent of Florida voters describe themselves as Hispanic, and about black. In a presidential-election battleground state that decided the 2000 race by gi Bush a margin of only 537 votes, the effect could be significant: black voters a Democratic, while Hispanics in Florida tend to vote Republican. Elections officials of Florida's Republican administration denied any partisan method they adopted, and noted that it had been approved as part of a settler lawsuit. 	PRINT PRINT SINGLE SAVE SAVE ARTICLE TT SHARE ARTICLE TT ARTICLE T	nt as rge W. helmingly use of the civil rights	F 9 9 9 9	Published (54 & up ((269 <u>Late</u> (199 & up (313 & up (327 & up (123 & up (169/mt <u>At</u>	on 08.08.07 Systemwide summer Lo Biggest cru Last-minute Affordable Weekend fl hens & Gree	US fare os Cabos iise sale e flights Labor Da lights to ek island	sale getaw of the s to Europ av vacat Vegas I combo





The controversial Terrorism Information Awareness program was conceived by retired Adm. John Poindexter and was run by the Information Awareness Office that he headed inside the Defense Advanced Research Projects Agency. It was developing software that could examine the computerized travel, credit, medical and other records of Americans and others around the world to search for telltale activities that might reveal preparations for a terrorist attack.

citizens aboard or foreigners in this country.

Sen. Ron Wyden (D-Oregon), who has led a campaign against the program, hailed the result Wednesday. "Americans on American soil are not going to be targets of TIA surveillance that would have violated their privacy and civil liberties. The government is not going to be able to pick Americans up by their ankles and shake them to see if anything funny falls out," Wyden said in an interview.







FINAL REPORT OF THE NATIONAL COMMISSION ON TERRORIST ATTACKS UPON THE UNITED STATES



Section 13.3: UNITY OF EFFORT IN SHARING INFORMATION

The U.S. government has access to a vast amount of information. When databases not usually thought of as "intelligence," such as customs or immigration information, are included, the storehouse is immense. ... In interviews around the government, official after official urged us to call attention to frustrations with the unglamorous "back office" side of government operations. ...

Recommendation: The president should lead the government-wide effort to bring the major national security institutions into the information revolution. He should coordinate the resolution of the legal, policy, and technical issues across agencies to create a "trusted information network."

AUTHORIZED EDITION

• Socio-Political

- Socio-Political
 - Bad data mining

- Socio-Political
 - Bad data mining
 - Misused data mining

- Socio-Political
 - Bad data mining
 - Misused data mining
 - Ignorant decision-making

- Socio-Political
 - Bad data mining
 - Misused data mining
 - Ignorant decision-making
 - Ramifications of data mining

- Socio-Political
 - Bad data mining
 - Misused data mining
 - Ignorant decision-making
 - Ramifications of data mining
 - Presuming fixed technology

- Data Mining is about Real Data: Benchmark data sets are a means to an end
 - Data sets are supposed to be representative of the sorts of problems our algorithms will see in practice
 - Data sets must stay timely as technological and scientific advances allow our ambitions to grow
 - A data set from some domain is not an application Who do you personally know that cares about your results?

- How do we ensure reproducible results?
 - Many of the applications of data mining are in the commercial sector -- How do we handle research results that reflect proprietary or otherwise restricted data?
 - How do we make sure academic research results address problems that are important in practice?
 - How do we handle inherent resource differentials between industry and academic research?
 - Access to data
 - Massive data centers
 - What new models of publication are particularly suited to data mining – "Executable articles" (Mark Liberman)