CRS Report for Congress

Received through the CRS Web

Data Mining: An Overview

Updated December 16, 2004

Jeffrey W. Seifert Analyst in Information Science and Technology Policy Resources, Science, and Industry Division

Data Mining: An Overview

Summary

Data mining is emerging as one of the key features of many homeland security initiatives. Often used as a means for detecting fraud, assessing risk, and product retailing, data mining involves the use of data analysis tools to discover previously unknown, valid patterns and relationships in large data sets. In the context of homeland security, data mining is often viewed as a potential means to identify terrorist activities, such as money transfers and communications, and to identify and track individual terrorists themselves, such as through travel and immigration records.

While data mining represents a significant advance in the type of analytical tools currently available, there are limitations to its capability. One limitation is that although data mining can help reveal patterns and relationships, it does not tell the user the value or significance of these patterns. These types of determinations must be made by the user. A second limitation is that while data mining can identify connections between behaviors and/or variables, it does not necessarily identify a causal relationship. To be successful, data mining still requires skilled technical and analytical specialists who can structure the analysis and interpret the output that is created.

Data mining is becoming increasingly common in both the private and public sectors. Industries such as banking, insurance, medicine, and retailing commonly use data mining to reduce costs, enhance research, and increase sales. In the public sector, data mining applications initially were used as a means to detect fraud and waste, but have grown to also be used for purposes such as measuring and improving program performance. However, some of the homeland security data mining applications represent a significant expansion in the quantity and scope of data to be analyzed. Two efforts that have attracted a higher level of congressional interest include the Terrorism Information Awareness (TIA) project (now-discontinued) and the Computer-Assisted Passenger Prescreening System II (CAPPS II) project (now-canceled and replaced by Secure Flight).

As with other aspects of data mining, while technological capabilities are important, there are other implementation and oversight issues that can influence the success of a project's outcome. One issue is data quality, which refers to the accuracy and completeness of the data being analyzed. A second issue is the interoperability of the data mining software and databases being used by different agencies. A third issue is mission creep, or the use of data for purposes other than for which the data were originally collected. A fourth issue is privacy. Questions that may be considered include the degree to which government agencies should use and mix commercial data with government data, whether data sources are being used for purposes other than those for which they were originally designed, and possible application of the Privacy Act to these initiatives. It is anticipated that congressional oversight of data mining projects will grow as data mining efforts continue to evolve. This report will be updated as events warrant.

Contents

What is Data Mining?	. 1
Limitations of Data Mining	. 3
Data Mining Uses	. 3
Terrorism Information Awareness (TIA) Program	
Computer-Assisted Passenger Prescreening System (CAPPS II)	
Data Mining Issues	11
Data Quality	
Interoperability	
Mission Creep	
Privacy	13
Legislation in the 108 th Congress	13
For Further Reading	16

Data Mining: An Overview

What is Data Mining?

Data mining involves the use of sophisticated data analysis tools to discover previously unknown, valid patterns and relationships in large data sets.¹ These tools can include statistical models, mathematical algorithms, and machine learning methods (algorithms that improve their performance automatically through experience, such as neural networks or decision trees). Consequently, data mining consists of more than collecting and managing data, it also includes analysis and prediction.

Data mining can be performed on data represented in quantitative, textual, or multimedia forms. Data mining applications can use a variety of parameters to examine the data. They include association (patterns where one event is connected to another event, such as purchasing a pen and purchasing paper), sequence or path analysis (patterns where one event leads to another event, such as the birth of a child and purchasing diapers), classification (identification of new patterns, such as coincidences between duct tape purchases and plastic sheeting purchases), clustering (finding and visually documenting groups of previously unknown facts, such as geographic location and brand preferences), and forecasting (discovering patterns from which one can make reasonable predictions regarding future activities, such as the prediction that people who join an athletic club may take exercise classes).²

As an application, compared to other data analysis applications, such as structured queries (used in many commercial databases) or statistical analysis software, data mining represents a *difference of kind rather than degree*. Many simpler analytical tools utilize a verification-based approach, where the user develops a hypothesis and then tests the data to prove or disprove the hypothesis. For example, a user might hypothesize that a customer who buys a hammer, will also buy a box of nails. The effectiveness of this approach can be limited by the creativity of the user to develop various hypotheses, as well as the structure of the software being used. In contrast, data mining utilizes a discovery approach, in which algorithms can be used to examine several multidimensional data relationships simultaneously, identifying those that are unique or frequently represented. For example, a hardware store may compare their customers' tool purchases with home ownership, type of automobile driven, age, occupation, income, and/or distance between residence and

¹ Two Crows Corporation, *Introduction to Data Mining and Knowledge Discovery, Third Edition* (Potomac, MD: Two Crows Corporation, 1999); Pieter Adriaans and Dolf Zantinge, *Data Mining* (New York: Addison Wesley, 1996).

For a more technically-oriented definition of data mining, see [http://searchcrm.techtarget.com/gDefinition/0,294236,sid11_gci211901,00.html].

the store. As a result of its complex capabilities, two precursors are important for a successful data mining exercise; a clear formulation of the problem to be solved, and access to the relevant data.³

Reflecting this conceptualization of data mining, some observers consider data mining to be just one step in a larger process known as knowledge discovery in databases (KDD). Other steps in the KDD process, in progressive order, include data cleaning, data integration, data selection, data transformation, (data mining), pattern evaluation, and knowledge presentation.⁴

A number of advances in technology and business processes have contributed to a growing interest in data mining in both the public and private sectors. Some of these changes include the growth of computer networks, which can be used to connect databases; the development of enhanced search-related techniques such as neural networks and advanced algorithms; the spread of the client/server computing model, allowing users to access centralized data resources from the desktop; and an increased ability to combine data from disparate sources into a single searchable source.⁵

In addition to these improved data management tools, the increased availability of information and the decreasing costs of storing it have also played a role. Over the past several years there has been a rapid increase in the volume of information collected and stored, with some observers suggesting that the quantity of the world's data approximately doubles every year. At the same time, the costs of data storage have decreased significantly from dollars per megabyte to pennies per megabyte. Similarly, computing power has continued to double every 18-24 months, while the relative cost of computing power has continued to decrease.

Data mining has become increasingly common in both the public and private sectors. Organizations use data mining as a tool to survey customer information, reduce fraud and waste, and assist in medical research. However, the proliferation of data mining has raised some implementation and oversight issues as well. These include concerns about the quality of the data being analyzed, the interoperability of the databases and software between agencies, and potential infringements on privacy. Also, there are some concerns that the limitations of data mining are being overlooked as agencies work to emphasize their homeland security initiatives.

³ John Makulowich, "Government Data Mining Systems Defy Definition," *Washington Technology*, 22 February 1999, [http://www.washingtontechnology.com/news/13_22/tech_features/393-3.html].

⁴ Jiawei Han and Micheline Kamber, *Data Mining: Concepts and Techniques* (New York: Morgan Kaufmann Publishers, 2001), p. 7.

⁵ Pieter Adriaans and Dolf Zantinge, *Data Mining* (New York: Addison Wesley, 1996), pp. 5-6.

⁶ Ibid., p. 2.

⁷ Two Crows Corporation, *Introduction to Data Mining and Knowledge Discovery, Third Edition* (Potomac, MD: Two Crows Corporation, 1999), p. 4.

Limitations of Data Mining

While data mining products can be very powerful tools, they are not self-sufficient applications. To be successful, data mining requires skilled technical and analytical specialists who can structure the analysis and interpret the output that is created. Consequently, the limitations of data mining are primarily data or personnel-related, rather than technology-related.⁸

Although data mining can help reveal patterns and relationships, it does not tell the user the value or significance of these patterns. These types of determinations must be made by the user. Similarly, the validity of the patterns discovered is dependent on how they compare to "real world" circumstances. For example, to assess the validity of a data mining application designed to identify potential terrorist suspects in a large pool of individuals, the user may test the model using data that includes information about known terrorists. However, while possibly re-affirming a particular profile, it does not necessarily mean that the application will identify a suspect whose behavior significantly deviates from the original model.

Another limitation of data mining is that while it can identify connections between behaviors and/or variables, it does not necessarily identify a causal relationship. For example, an application may identify that a pattern of behavior, such as the propensity to purchase airline tickets just shortly before the flight is scheduled to depart, is related to characteristics such as income, level of education, and Internet use. However, that does not necessarily indicate that the ticket purchasing behavior is caused by one or more of these variables. In fact, the individual's behavior could be affected by some additional variable(s) such as occupation (the need to make trips on short notice), family status (a sick relative needing care), or a hobby (taking advantage of last minute discounts to visit new destinations).

Data Mining Uses

Data mining is used for a variety of purposes in both the private and public sectors. Industries such as banking, insurance, medicine, and retailing commonly use data mining to reduce costs, enhance research, and increase sales. For example, the insurance and banking industries use data mining applications to detect fraud and assist in risk assessment (e.g., credit scoring). Using customer data collected over several years, companies can develop models that predict whether a customer is a good credit risk, or whether an accident claim may be fraudulent and should be investigated more closely. The medical community sometimes uses data mining to help predict the effectiveness of a procedure or medicine. Pharmaceutical firms use data mining of chemical compounds and genetic material to help guide research on new treatments for diseases. Retailers can use information collected through affinity programs (e.g., shoppers' club cards, frequent flyer points, contests) to assess the

⁸ Ibid., p. 2.

⁹ Ibid., p. 1.

effectiveness of product selection and placement decisions, coupon offers, and which products are often purchased together. Companies such as telephone service providers and music clubs can use data mining to create a "churn analysis," to assess which customers are likely to remain as subscribers and which ones are likely to switch to a competitor.¹⁰

In the public sector, data mining applications were initially used as a means to detect fraud and waste, but they have grown also to be used for purposes such as measuring and improving program performance. It has been reported that data mining has helped the federal government recover millions of dollars in fraudulent Medicare payments.¹¹ The Justice Department has been able to use data mining to assess crime patterns and adjust resource allotments accordingly. Similarly, the Department of Veterans Affairs has used data mining to help predict demographic changes in the constituency it serves so that it can better estimate its budgetary needs. Another example is the Federal Aviation Administration, which uses data mining to review plane crash data to recognize common defects and recommend precautionary measures.¹²

Recently, data mining has been increasingly cited as an important tool for homeland security efforts. Some observers suggest that data mining should be used as a means to identify terrorist activities, such as money transfers and communications, and to identify and track individual terrorists themselves, such as through travel and immigration records. Two initiatives that have attracted significant attention include the now-discontinued Terrorism Information Awareness (TIA) project¹³ conducted by the Defense Advanced Research Projects Agency (DARPA), and the now-canceled Computer-Assisted Passenger Prescreening System II (CAPPS II) that was being developed by the Transportation Security Administration (TSA). CAPPS II is being replaced by a new program called Secure Flight.

Section 8131 of the FY2004 Department of Defense Appropriations Act (P.L. 108-87) prohibited further funding of TIA as a whole, while allowing unspecified subcomponents of the TIA initiative to be funded as part of DOD's classified budget, subject to the provisions of the National Foreign Intelligence Program, which restricts the processing and analysis of information on U.S. citizens. For further details regarding this provision, see CRS Report RL31805 *Authorization and Appropriations for FY2004: Defense*, by Amy Belasco and Stephen Daggett.

¹⁰ Two Crows Corporation, *Introduction to Data Mining and Knowledge Discovery, Third Edition* (Potomac, MD: Two Crows Corporation, 1999), p. 5; Patrick Dillon, *Data Mining: Transforming Business Data Into Competitive Advantage and Intellectual Capital* (Atlanta GA: The Information Management Forum, 1998), pp. 5-6.

¹¹ George Cahlink, "Data Mining Taps the Trends," *Government Executive Magazine*, October 1, 2000, [http://www.govexec.com/tech/articles/1000managetech.htm].

¹² Ibid.; for a more detailed review of the purpose for data mining conducted by federal departments and agencies, see U.S. General Accounting Office, *Data Mining: Federal Efforts Cover a Wide Range of Uses*, GAO Report GAO-04-548 (Washington: May 2004).

¹³ This project was originally identified as the Total Information Awareness project until DARPA publicly renamed it the Terrorism Information Awareness project in May 2003.

Terrorism Information Awareness (TIA) Program

In the immediate aftermath of the September 11, 2001, terrorist attacks, many questions were raised about the country's intelligence tools and capabilities, as well as the government's ability to detect other so-called "sleeper cells," if, indeed, they existed. One response to these concerns was the creation of the Information Awareness Office (IAO) at the Defense Advanced Research Projects Agency (DARPA)¹⁴ in January 2002. The role of IAO was "in part to bring together, under the leadership of one technical office director, several existing DARPA programs focused on applying information technology to combat terrorist threats."15 The mission statement for IAO suggested that the emphasis on these technology programs was to "counter asymmetric threats by achieving total information awareness useful for preemption, national security warning, and national security decision making."¹⁶ To that end, the TIA project was to focus on three specific areas of research, anticipated to be conducted over five years, to develop technologies that would assist in the detection of terrorist groups planning attacks against American interests, both inside and outside the country. The three areas of research and their purposes were described in a DOD Inspector General report as:

"... language translation, data search with pattern recognition and privacy protection, and advanced collaborative and decision support tools. Language translation technology would enable the rapid analysis of foreign languages, both spoken and written, and allow analysts to quickly search the translated materials for clues about emerging threats. The data search, pattern recognition, and privacy protection technologies would permit analysts to search vast quantities of data for patterns that suggest terrorist activity while at the same time controlling access to the data, enforcing laws and policies, and ensuring detection of misuse of the information obtained. The collaborative reasoning and decision support technologies would allow analysts from different agencies to share data."¹⁷

Each part had the potential to improve the data mining capabilities of agencies that adopt the technology.¹⁸ Automated rapid language translation could allow

¹⁴ DARPA "is the central research and development organization for the Department of Defense (DOD)" that engages in basic and applied research, with a specific focus on "research and technology where risk and payoff are both very high and where success may provide dramatic advances for traditional military roles and missions." [http://www.darpa.mil/]

¹⁵ Department of Defense. 20 May 2003. Report to Congress Regarding the Terrorism Information Awareness Program, Executive Summary, p.2.

¹⁶ Department of Defense. 20 May 2003. *Report to Congress Regarding the Terrorism Information Awareness Program, Detailed Information*, p.1 (emphasis added).

Department of Defense, Office of the Inspector General. 12 December 2003. *Information Technology Management: Terrorism Information Awareness Project (D2004-033)*. P. 7.

¹⁸ It is important to note that while DARPA's mission is to conduct research and development on technologies that can be used to address national-level problems, it would (continued...)

analysts to search and monitor foreign language documents and transmissions more quickly than currently possible. Improved search and pattern recognition technologies may enable more comprehensive and thorough mining of transactional data, such as passport and visa applications, car rentals, driver license renewals, criminal records, and airline ticket purchases. Improved collaboration and decision support tools might facilitate the search and coordination activities being conducted by different agencies and levels of government.¹⁹

In public statements DARPA frequently referred to the TIA program as a research and development project designed to create experimental prototype tools, and that the research agency would only use "data that is legally available and obtainable by the U.S. Government." DARPA further emphasized that these tools could be adopted and used by *other* agencies, and that DARPA itself would not be engaging in any actual-use data mining applications, although it could "support production of a scalable leave-behind system prototype." In addition, some of the technology projects being carried out in association with the TIA program did not involve data mining. However, the TIA program's overall emphasis on collecting, tracking, and analyzing data trails left by individuals served to generate significant and vocal opposition soon after John Poindexter made a presentation on TIA at the DARPATech 2002 Conference in August 2002.²³

Critics of the TIA program were further incensed by two administrative aspects of the project. The first involved the Director of IAO, Dr. John M. Poindexter. Poindexter, a retired Admiral, was, until that time, perhaps most well-known for his alleged role in the Iran-contra scandal during the Reagan Administration. His involvement with the program caused many in the civil liberties community to

¹⁸ (...continued) not be responsible for the operation of TIA, if it were to be adopted.

¹⁹ For more details about the Terrorism Information Awareness program and related information and privacy laws, see CRS Report RL31730, *Privacy: Total Information Awareness Programs and Related Information Access, Collection, and Protection Laws*, by Gina Marie Stevens, and CRS Report RL31786, *Total Information Awareness Programs: Funding, Composition, and Oversight Issues*, by Amy Belasco.

Department of Defense, DARPA, "Defense Advanced Research Project Agency's Information Awareness Office and Total Information Awareness Project," p. 1, [http://www.iwar.org.uk/news-archive/tia/iaotia.pdf].

²¹ Ibid., p. 2.

²² Although most of the TIA-related projects did involve some form of data collection, the primary purposes of some of these projects, such as war gaming, language translation, and biological agent detection, were less connected to data mining activities. For a description of these projects, see [http://www.fas.org/irp/agency/dod/poindexter.html].

The text of Poindexter's presentation is available at [http://www.darpa.mil/DARPATech2002/presentations/iao_pdf/speeches/POINDEXT.pdf]. The slide presentation of Poindexter's presentation is available at [http://www.darpa.mil/DARPATech2002/presentations/iao_pdf/slides/PoindexterIAO.pdf].

question the true motives behind TIA.²⁴ The second source of contention involved TIA's original logo, which depicted an "all-seeing" eye atop of a pyramid looking down over the globe, accompanied by the Latin phrase *scientia est potentia* (knowledge is power).²⁵ Although DARPA eventually removed the logo from its website, it left a lasting impression.

The continued negative publicity surrounding the TIA program contributed to the introduction of a number of bills in Congress that eventually led to the program's dissolution. Among these bills was S. 188, the Data-Mining Moratorium Act of 2003, which, if passed, would have imposed a moratorium on the implementation of data mining under the TIA program by the Department of Defense, as well as any similar program by the Department of Homeland Security. An amendment included in the Omnibus Appropriations Act for Fiscal Year 2003 (P.L. 108-7) required the Director of Central Intelligence, the Secretary of Defense, and the Attorney General to submit a joint report to Congress within 90 days providing details about the TIA program.²⁶ Funding for TIA as a whole was prohibited with the passage of the FY2004 Department of Defense Appropriations Act (P.L. 108-87) in September 2003. However, Section 8131 of the law allowed unspecified subcomponents of the TIA initiative to be funded as part of DOD's classified budget, subject to the provisions of the National Foreign Intelligence Program, which restricts the processing and analysis of information on U.S. citizens.²⁷

Computer-Assisted Passenger Prescreening System (CAPPS II)

Similar to TIA, the CAPPS II project represented a direct response to the September 11, 2001, terrorist attacks. With the images of airliners flying into buildings fresh in people's minds, air travel was now widely viewed not only as a critically vulnerable terrorist target, but also as a weapon for inflicting larger harm. The CAPPS II initiative was intended to replace the original CAPPS, currently being used. Spurred, in part, by the growing number of airplane bombings, the existing CAPPS (originally called CAPS) was developed through a grant provided by the Federal Aviation Administration (FAA) to Northwest Airlines, with a prototype

²⁴ Shane Harris, "Counterterrorism Project Assailed By Lawmakers, Privacy Advocates," *Government Executive Magazine*, 25 November 2002, [http://www.govexec.com/dailyfed/1102/112502h1.htm].

²⁵ The original logo can be found at [http://www.thememoryhole.org/policestate/iaologo.htm].

²⁶ The report is available at [http://www.eff.org/Privacy/TIA/TIA-report.pdf]. Some of the information required includes spending schedules, likely effectiveness of the program, likely impact on privacy and civil liberties, and any laws and regulations that may need to be changed to fully deploy TIA. If the report had not submitted within 90 days, funding for the TIA program could have been discontinued. For more details regarding this amendment, see CRS Report RL31786, *Total Information Awareness Programs: Funding, Composition, and Oversight Issues*, by Amy Belasco.

²⁷ For further details regarding this provision, see CRS Report RL31805 *Authorization and Appropriations for FY2004: Defense*, by Amy Belasco and Stephen Daggett.

systems, and, by 1998, most of the U.S.-based airlines had voluntarily implemented CAPS, with the remaining few working toward implementation. Also, during this time, the White House Commission on Aviation Safety and Security (sometimes referred to as the Gore Commission) released its final report in February 1997. Included in the commission's report was a recommendation that the United States implement automated passenger profiling for its airports. On April 19, 1999, the FAA issued a notice of proposed rulemaking (NPRM) regarding the security of checked baggage on flights within the United States (docket no. FAA-1999-5536). As part of this still-pending rule, domestic flights would be required to utilize "the FAA-approved computer-assisted passenger screening (CAPS) system to select passengers whose checked baggage must be subjected to additional security measures."

The current CAPPS system is a rule-based system that uses the information provided by the passenger when purchasing the ticket to determine if the passenger fits into one of two categories; "selectees" requiring additional security screening, and those who do not. CAPPS also compares the passenger name to those on a list of known or suspected terrorists. CAPPS II was described by TSA as "an enhanced system to confirm the identities of passengers and to identify foreign terrorists or persons with terrorist connections before they can board U.S. aircraft. CAPPS II would have sent information provided by the passenger in the passengers name record (PNR), including full name, address, phone number, and date of birth, to commercial data providers for comparison to authenticate the identity of the passenger. The commercial data provider would have then transmitted a numerical score back to TSA indicating a particular risk level. Passengers with a "green" score would have undergone "normal screening," while passengers with a "green" score would have undergone additional screening. Passengers with a "red" score would not have been allowed to board the flight, and would have received "the

Department of Transportation, *White House Commission on Aviation and Security: The DOT Status Report*, February 1998, [http://www.dot.gov/affairs/whcoasas.htm].

²⁹ The Gore Commission was established by Executive Order 13015 on August 22, 1996, following the crash of TWA flight 800 in July 1996.

White House Commission on Aviation Safety and Security: Final Report to President Clinton. 12 February 1997. [http://www.fas.org/irp/threat/212fin~1.html].

The docket can be found online at [http://dms.dot.gov/search/document.cfm?documentid=57279&docketid=5536].

³² Federal Register, 64 (April 19,1999): 19220.

³³ U.S. General Accounting Office, *Aviation Security: Computer-Assisted Passenger Prescreening System Faces Significant Implementation Challenges*, GAO Report GAO-04-385, February 2004, pp. 5-6.

Transportation Security Administration, "TSA's CAPPS II Gives Equal Weight to Privacy, Security," Press Release, 11 March 2003, [http://www.tsa.gov/public/display?theme=44&content=535].

³⁵ Robert O'Harrow, Jr., "Aviation ID System Stirs Doubt," *Washington Post*, 14 March 2003, p. A16.

attention of law enforcement."³⁶ While drawing on information from commercial databases, TSA had stated that it would not see the actual information used to calculate the scores, and that it would not retain the traveler's information.

TSA had planned to test the system at selected airports during spring 2004.³⁷ However, CAPPS II encountered a number of obstacles to implementation. One obstacle involved obtaining the required data to test the system. Several high-profile debacles resulting in class-action lawsuits have made the U.S.-based airlines very wary of voluntarily providing passenger information. In early 2003, Delta Airlines was to begin testing CAPPS II using its customers' passenger data at three airports across the country. However, Delta became the target of a vociferous boycott campaign, raising further concerns about CAPPS II generally. 38 In September 2003, it was revealed that JetBlue shared private passenger information in September 2002 with Torch Concepts, a defense contractor, which was testing a data mining application for the U.S. Army. The information shared reportedly included itineraries, names, addresses, and phone numbers for 1.5 million passengers.³⁹ In January 2004, it was reported that Northwest Airlines provided personal information on millions of its passengers to the National Aeronautics and Space Administration (NASA) from October to December 2001 for an airline security-related data mining experiment.⁴⁰ In April 2004, it was revealed that American Airlines agreed to provide private passenger data on 1.2 million of its customers to TSA in June 2002, although the information was sent instead to four companies competing to win a contract with TSA.⁴¹ Further instances of data being provided for the purpose of testing CAPPS II were brought to light during a Senate Committee on Government Affairs confirmation hearing on June 23, 2004. In his answers to the committee, the acting director of TSA, David M. Stone, stated that during 2002 and 2003 four airlines; Delta, Continental, America West, and Frontier, and two travel reservation companies; Galileo International and Sabre Holdings, provided passenger records to TSA and/or its contractors.⁴²

Transportation Security Administration, "TSA's CAPPS II Gives Equal Weight to Privacy, Security," Press Release, 11 March 2003, [http://www.tsa.gov/public/display?theme=44&content=535].

³⁷ Sara Kehaulani Goo, "U.S. to Push Airlines for Passenger Records," *Washington Post*, 12 January 2004, p. A1.

The Boycott Delta website is available at [http://www.boycottdelta.org].

³⁹ Don Phillips, "JetBlue Apologizes for Use of Passenger Records," *The Washington Post*, 20 September 2003, p. E1; Sara Kehaulani Goo, "TSA Helped JetBlue Share Data, Report Says," *Washington Post*, 21 February 2004, p. E1.

⁴⁰ Sara Kehaulani Goo, "Northwest Gave U.S. Data on Passengers," *Washington Post*,18 January 2004, p. A1.

⁴¹ Sara Kehaulani Goo, "American Airlines Revealed Passenger Data," *Washington Post*, 10 April 2004, p. D12.

For the written responses to the committee's questions, see [http://www.epic.org/privacy/airtravel/stone_answers.pdf]; Sara Kehaulani Goo, "Agency Got More Airline Records," *Washington Post*, 24 June 2004, p. A16.

Concerns about privacy protections had also dissuaded the European Union (EU) from providing any data to TSA to test CAPPS II. However, in May 2004, the EU signed an agreement with the United States that would have allowed PNR data for flights originating from the EU to be used in testing CAPPS II, but only after TSA was authorized to use domestic data as well. As part of the agreement, the EU data was to be retained for only three-and-a-half years (unless it is part of a law enforcement action), only 34 of the 39 elements of the PNR were to be accessed by authorities, and there were to be yearly joint DHS-EU reviews of the implementation of the agreement.

Another obstacle was the perception of mission creep. CAPPS II was originally intended to just screen for high-risk passengers who may pose a threat to safe air travel. However, in an August 1, 2003, *Federal Register* notice, TSA stated that CAPPS II could also be used to identify individuals with outstanding state or federal arrest warrants, as well as identify both foreign *and* domestic terrorists (not just foreign terrorists). The notice also states that CAPPS II could be "linked with the U.S. Visitor and Immigrant Status Indicator Technology (US-VISIT) program" to identify individuals who are in the country illegally (e.g., individuals with expired visas, illegal aliens, etc.).⁴⁵ In response to critics who cited these possible uses as examples of mission creep, TSA claimed that the suggested uses were consistent with the goals of improving aviation security.⁴⁶

Several other concerns had also been raised, including the length of time passenger information was to be retained, who would have access to the information, the accuracy of the commercial data being used to authenticate a passenger's identity, the creation of procedures to allow passengers the opportunity to correct data errors in their records, and the ability of the system to detect attempts by individuals to use identity theft to board a plane undetected.

In August 2004, TSA announced that the CAPPS II program was being canceled and would be replaced with a new system called Secure Flight. In the Department of Homeland Security Appropriations Act, 2005 (P.L. 108-334), Congress included a provision (Sec. 522) prohibiting the use of appropriated funds for "deployment or implementation, on other than a test basis," of CAPPS II, Secure Flight, "or other follow on/successor programs," until GAO has certified that such a system has met

⁴³ Some information, such as meal preferences, which could be used to infer religious affiliation, and health considerations will not be made available. Goo, Sara Kehaulani, "U.S., EU Will Share Passenger Records," *Washington Post*, 29 May 2004, p. A2.

Department of Homeland Security, "Fact Sheet: US-EU Passenger Name Record Agreement Signed," 28 May 2004, [http://www.dhs.gov/dhspublic/display?content=3651].

⁴⁵ Federal Register. Vol. 68 No. 148 Friday August 1, 2003. P. 45266; U.S. General Accounting Office, Aviation Security: Challenges Delay Implementation of Computer-Assisted Passenger Prescreening System, GAO Testimony GAO-04-504T, 17 March 2004, p. 17

⁴⁶ U.S. General Accounting Office, Aviation Security: Challenges Delay Implementation of Computer-Assisted Passenger Prescreening System, GAO Testimony GAO-04-504T, 17 March 2004, p. 17

all eight of the privacy requirements enumerated in a February 2004 GAO report,⁴⁷ can accommodate any unique air transportation needs as it relates to interstate transportation, and that "appropriate life-cycle cost estimates, and expenditure and program plans exist." GAO's certification report is due to Congress no later than March 28, 2005.

Data Mining Issues

As data mining initiatives continue to evolve, there are several issues Congress may decide to consider related to implementation and oversight. These issues include, but are not limited to, data quality, interoperability, mission creep, and privacy. As with other aspects of data mining, while technological capabilities are important, other factors also influence the success of a project's outcome.

Data Quality

Data quality is a multifaceted issue that represents one of the biggest challenges for data mining. Data quality refers to the accuracy and completeness of the data. Data quality can also be affected by the structure and consistency of the data being analyzed. The presence of duplicate records, the lack of data standards, the timeliness of updates, and human error can significantly impact the effectiveness of the more complex data mining techniques, which are sensitive to subtle differences that may exist in the data. To improve data quality, it is sometimes necessary to "clean" the data, which can involve the removal of duplicate records, normalizing the values used to represent information in the database (e.g., ensuring that "no" is represented as a 0 throughout the database, and not sometimes as a 0, sometimes as a N, etc.), accounting for missing data points, removing unneeded data fields, identifying anomalous data points (e.g., an individual whose age is shown as 142 years), and standardizing data formats (e.g., changing dates so they all include MM/DD/YYYY).

Interoperability

Related to data quality, is the issue of interoperability of different databases and data mining software. Interoperability refers to the ability of a computer system and/or data to work with other systems or data using common standards or processes. Interoperability is a critical part of the larger efforts to improve interagency collaboration and information sharing through e-government and homeland security initiatives. For data mining, interoperability of databases and software is important to enable the search and analysis of multiple databases simultaneously, and to help ensure the compatibility of data mining activities of different agencies. Data mining projects that are trying to take advantage of existing legacy databases or that are

⁴⁷ The eight issues included establishing an oversight board, ensuring the accuracy of the data used, conducting stress testing, instituting abuse prevention practices, preventing unauthorized access, establishing clear policies for the operation and use of the system, satisfying privacy concerns, and created a redress process. U.S. General Accounting Office, *Aviation Security: Computer-Assisted Passenger Prescreening System Faces Significant Implementation Challenges*, GAO Report GAO-04-385, February 2004.

initiating first-time collaborative efforts with other agencies or levels of government (e.g., police departments in different states) may experience interoperability problems. Similarly, as agencies move forward with the creation of new databases and information sharing efforts, they will need to address interoperability issues during their planning stages to better ensure the effectiveness of their data mining projects.

Mission Creep

Mission creep is one of the leading risks of data mining cited by civil libertarians, and represents how control over one's information can be a tenuous proposition. Mission creep refers to the use of data for purposes other than that for which the data was originally collected. This can occur regardless of whether the data was provided voluntarily by the individual or was collected through other means.

Efforts to fight terrorism can, at times, take on an acute sense of urgency. This urgency can create pressure on both data holders and officials who access the data. To leave an available resource unused may appear to some as being negligent. Data holders may feel obligated to make any information available that could be used to prevent a future attack or track a known terrorist. Similarly, government officials responsible for ensuring the safety of others may be pressured to use and/or combine existing databases to identify potential threats. Unlike physical searches, or the detention of individuals, accessing information for purposes other than originally intended may appear to be a victimless or harmless exercise. However, such information use can lead to unintended outcomes and produce misleading results.

One of the primary reasons for misleading results is inaccurate data. All data collection efforts suffer accuracy concerns to some degree. Ensuring the accuracy of information can require costly protocols that may not be cost effective if the data is not of inherently high economic value. In well-managed data mining projects, the original data collecting organization is likely to be aware of the data's limitations and account for these limitations accordingly. However, such awareness may not be communicated or heeded when data is used for other purposes. For example, the accuracy of information collected through a shopper's club card may suffer for a variety of reasons, including the lack of identity authentication when a card is issued, cashiers using their own cards for customers who do not have one, and/or customers who use multiple cards. ⁴⁸ For the purposes of marketing to consumers, the impact of these inaccuracies is negligible to the individual. If a government agency were to use that information to target individuals based on food purchases associated with particular religious observances though, an outcome based on inaccurate information could be, at the least, a waste of resources by the government agency, and an unpleasant experience for the misidentified individual. As the March 2004 TAPAC report observes, the potential wide reuse of data suggests that concerns about mission creep can extend beyond privacy to the protection of civil rights in the event that information is used for "targeting an individual solely on the basis of religion or

⁴⁸ Technology and Privacy Advisory Committee, Department of Defense. *Safeguarding Privacy in the Fight Against Terrorism*, March 2004, p. 40.

expression, or using information in a way that would violate the constitutional guarantee against self-incrimination."⁴⁹

Privacy

As additional information sharing and data mining initiatives have been announced, increased attention has focused on the implications for privacy. Concerns about privacy focus both on actual projects proposed, as well as concerns about the potential for data mining applications to be expanded beyond their original purposes (mission creep). For example, some experts suggest that anti-terrorism data mining applications might also be useful for combating other types of crime as well.⁵⁰ So far there has been little consensus about how data mining should be carried out, with several competing points of view being debated. Some observers contend that tradeoffs may need to be made regarding privacy to ensure security. Other observers suggest that existing laws and regulations regarding privacy protections are adequate, and that these initiatives do not pose any threats to privacy. Still other observers argue that not enough is known about how data mining projects will be carried out, and that greater oversight is needed. There is also some disagreement over how privacy concerns should be addressed. Some observers suggest that technical solutions are adequate. In contrast, some privacy advocates argue in favor of creating clearer policies and exercising stronger oversight. As data mining efforts move forward, Congress may consider a variety of questions including, the degree to which government agencies should use and mix commercial data with government data, whether data sources are being used for purposes other than those for which they were originally designed, and the possible application of the Privacy Act to these initiatives.

Legislation in the 108th Congress

During the 108th Congress, a number of legislative proposals were introduced that would restrict data mining activities by some parts of the federal government, and/or increase the reporting requirements of such projects to Congress. For example, on January 16, 2003, Senator Feingold introduced S. 188 the Data-Mining Moratorium Act of 2003, which would have imposed a moratorium on the implementation of data mining under the Total Information Awareness program (now referred to as the Terrorism Information Awareness project) by the Department of Defense, as well as any similar program by the Department of Homeland Security. S. 188 was referred to the Committee on the Judiciary.

On January 23, 2003, Senator Wyden introduced S.Amdt. 59, an amendment to H.J.Res. 2, the Omnibus Appropriations Act for Fiscal Year 2003. As passed in its final form as part of the omnibus spending bill (P.L. 108-7) on February 13, 2003,

⁴⁹ Ibid., p. 39.

⁵⁰ Drew Clark, "Privacy Experts Differ on Merits of Passenger-Screening Program," *Government Executive Magazine*, November 21, 2003, [http://www.govexec.com/dailyfed/1103/112103td2.htm].

and signed by the President on February 20, 2003, the amendment requires the Director of Central Intelligence, the Secretary of Defense, and the Attorney General to submit a joint report to Congress within 90 days providing details about the TIA program. Some of the information required includes spending schedules, likely effectiveness of the program, likely impact on privacy and civil liberties, and any laws and regulations that may need to be changed to fully deploy TIA. If the report had not submitted within 90 days, funding for the TIA program could have been discontinued. Funding for TIA was later discontinued in Section 8131 of the FY2004 Department of Defense Appropriations Act (P.L. 108-87), signed into law on September 30, 2003.

On March 13, 2003, Senator Wyden introduced an amendment to S. 165 the Air Cargo Security Act, requiring the Secretary of Homeland Security to submit a report to Congress within 90 days providing information about the impact of CAPPS II on privacy and civil liberties. The amendment was passed by the Committee on Commerce, Science, and Transportation, and the bill was forwarded for consideration by the full Senate (S.Rept. 108-38). In May 2003, S. 165 was passed by the Senate with the Wyden amendment included and was sent to the House where it was referred to the Committee on Transportation and Infrastructure.

Funding restrictions on CAPPSII were included in section 519 of the FY2004 Department of Homeland Security Appropriations Act (P.L. 108-90), signed into law October 1, 2003. This provision included restrictions on the "deployment or implementation, on other than a test basis, of the Computer-Assisted Passenger Prescreening System (CAPPSII)," pending the completion of a GAO report regarding the efficacy, accuracy, and security of CAPPSII, as well as the existence of a system of an appeals process for individuals identified as a potential threat by the system. ⁵⁴ In its report delivered to Congress in February 2004, GAO reported that "As of January 1, 2004, TSA has not fully addressed seven of the eight CAPPSII issues identified by the Congress as key areas of interest." The one issue GAO determined that TSA had addressed is the establishment of an internal oversight board. GAO

⁵¹ The report is available at [http://www.eff.org/Privacy/TIA/TIA-report.pdf].

⁵² For more details regarding this amendment, see CRS Report RL31786, *Total Information Awareness Programs: Funding, Composition, and Oversight Issues*, by Amy Belasco.

⁵³ For further details regarding this provision, see CRS Report RL31805 *Authorization and Appropriations for FY2004: Defense*, by Amy Belasco and Stephen Daggett.

⁵⁴ Section 519 of P.L. 108-90 specifically identifies eight issues that TSA must address before it can spend funds to deploy or implement CAPPSII on other than a test basis. These include 1. establishing a system of due process for passengers to correct erroneous information; 2. assess the accuracy of the databases being used; 3. stress test the system and demonstrate the efficiency and accuracy of the search tools; 4. establish and internal oversight board; 5. install operational safeguards to prevent abuse; 6. install security measures to protect against unauthorized access by hackers or other intruders; 7. establish policies for effective oversight of system use and operation; and 8. address any privacy concerns related to the system.

⁵⁵ General Accounting Office, Aviation Security: Computer-Assisted Passenger Prescreening System Faces Significant Implementation Challenges, GAO-04-385, February 2004, p. 4.

attributed the incomplete progress on these issues partly to the "early stage of the system's development." ⁵⁶

On March 25, 2003, the House Committee on Government Reform Subcommittee on Technology, Information Policy, Intergovernmental Relations, and the Census held a hearing on the current and future possibilities of data mining. The witnesses, drawn from federal and state government, industry, and academia, highlighted a number of perceived strengths and weaknesses of data mining, as well as the still-evolving nature of the technology and practices behind data mining. ⁵⁷ While data mining was alternatively described by some witnesses as a process, and by other witnesses as a productivity tool, there appeared to be a general consensus that the challenges facing the future development and success of government data mining applications were related less to technological concerns than to other issues such as data integrity, security, and privacy. On May 6 and May 20, 2003 the Subcommittee also held hearings on the potential opportunities and challenges for using factual data analysis for national security purposes.

On July 29, 2003 Senator Wyden introduced S. 1484 The Citizens' Protection in Federal Databases Act, which was referred to the Committee on the Judiciary. Among its provisions, S. 1484 would have required the Attorney General, the Secretary of Defense, the Secretary of Homeland Security, the Secretary of the Treasury, the Director of Central Intelligence, and the Director of the Federal Bureau of Investigation to submit to Congress a report containing information regarding the purposes, type of data, costs, contract durations, research methodologies, and other details before obligating or spending any funds on commercially available databases. S. 1484 would also have set restrictions on the conduct of searches or analysis of databases "based solely on a hypothetical scenario or hypothetical supposition of who may commit a crime or pose a threat to national security."

On July 31, 2003 Senator Feingold introduced S. 1544 the Data-Mining Reporting Act of 2003, which was referred to the Committee on the Judiciary. Among its provisions, S. 1544 would have required any department or agency engaged in data mining to submit a public report to Congress regarding these activities. These reports would have been required to include a variety of details about the data mining project, including a description of the technology and data to be used, an assessment of the expected efficacy of the data mining project, a privacy impact assessment, an analysis of the relevant laws and regulations that would govern the project, and a discussion of procedures for informing individuals their personal information will be used and allowing them to opt out, or an explanation of why such procedures are not in place.

Also on July 31, 2003, Senator Murkowski introduced S. 1552 the Protecting the Rights of Individuals Act, which was referred to the Committee on the Judiciary.

⁵⁶ Ibid

⁵⁷ Witnesses testifying at the hearing included Florida State Senator Paula Dockery, Dr. Jen Que Louie representing Nautilus Systems, Inc., Mark Forman representing OMB, Gregory Kutz representing GAO, and Jeffrey Rosen, an Associate Professor at George Washington University Law School.

Among its provisions, section 7 of S. 1552 would have imposed a moratorium on data mining by any federal department or agency "except pursuant to a law specifically authorizing such data-mining program or activity by such department or agency." It also would have required

The head of each department or agency of the Federal Government that engages or plans to engage in any activities relating to the development or use of a datamining program or activity shall submit to Congress, and make available to the public, a report on such activities.

On May 5, 2004, Representative McDermott introduced H.R. 4290 the Data-Mining Reporting Act of 2004, which was referred to the House Committee on Government Reform Subcommittee on Technology, Information Policy, Intergovernmental Relations, and the Census. H.R. 4290 would have required

each department or agency of the Federal Government that is engaged in any activity or use or develop data-mining technology shall each submit a public report to Congress on all such activities of the department or agency under the jurisdiction of that official.

A similar provision was included in H.R. 4591/S. 2528 the Civil Liberties Restoration Act of 2004. S. 2528 was introduced by Senator Kennedy on June 16, 2004 and referred to the Committee on the Judiciary. H.R. 4591 was introduced by Representative Berman on June 16, 2004 and referred to the Committee on the Judiciary and the Permanent Select Committee on Intelligence.

For Further Reading

- CRS Report RL32597, *Information Sharing for Homeland Security: A Brief Overview*, by Harold C. Relyea and Jeffrey W. Seifert.
- CRS Report RL31408, *Internet Privacy: Overview and Pending Legislation*, by Marcia S. Smith.
- CRS Report RL30671, *Personal Privacy Protection: The Legislative Response*, by Harold C. Relyea. Archived.
- CRS Report RL31730, *Privacy: Total Information Awareness Programs and Related Information Access, Collection, and Protection Laws*, by Gina Marie Stevens.
- CRS Report RL31786, Total Information Awareness Programs: Funding, Composition, and Oversight Issues, by Amy Belasco.
- DARPA, Report to Congress Regarding the Terrorism Information Awareness Program, May 20, 2003, [http://www.eff.org/Privacy/TIA/TIA-report.pdf].
- Department of Defense, Office of the Inspector General, *Information Technology Management: Terrorism Information Awareness Program (D-2004-033)*, December 12, 2003, [http://www.dodig.osd.mil/audit/reports/FY04/04-033.pdf].