

Incorporating Data Stream Analysis into Decision Support Systems

Damianos Chatziantoniou

Athens University of Economics and Business, Greece

George Doukidis

Athens University of Economics and Business, Greece

INTRODUCTION

Traditional decision support systems (DSS) and executive information systems (EIS) gather and present information from several sources for business purposes. It is an information technology to help the knowledge worker (executive, manager, analyst) make faster and better decisions. So far, this data was stored statically and persistently in a database, typically in a data warehouse. Data warehouses collect masses of operational data, allowing analysts to extract information by issuing decision support queries on the otherwise discarded data. In a typical scenario, an organization stores a detailed record of its operations in a database, which is then analyzed to improve efficiency, detect sales opportunities, and so on. Performing complex analysis on this data is an essential component of these organizations' businesses. Chaudhuri and Dayal (1997), present an excellent survey on decision-making and on-line analytical processing (OLAP) technologies for traditional database systems.

In many applications, however, it may not be possible to process queries within a database management system (DBMS). These applications involve data items that arrive on-line from multiple sources in a continuous, rapid and time-varying fashion (Babcock, Babu, Datar, Motwani & Widom, 2002). This data may or may not be stored in a database. As a result, a new class of data-intensive applications has recently attracted a lot of attention: applications in which the data is modeled not as persistent relations, but rather as transient *data streams*. Examples include financial applications (streams of transactions or ticks), network monitoring (stream of packets), security, telecommunication data management (stream of calls or call packets), web applications (clickstreams), manufacturing, sensor networks (measurements), and others. In data streams we usually have "continuous" queries (Terry, Goldberg, Nichols & Oki, 1992; Babu & Widom, 2001) rather than "one-time". The answer to a continuous query is produced over time, reflecting the stream data seen so far. Answers may be stored and updated as new data arrives or may be produced as data streams themselves.

Continuous queries can be used for monitoring, alerting, security, personalization, etc. Data streams can be either *transactional*, i.e. log interactions between entities, such as credit card purchases; web clickstreams; phone calls; or *measurement*, i.e. monitor evolution of entity states, such as physical phenomena, road traffic, temperature, network.

How to best model, express and evaluate complex queries over data streams is an open and difficult problem. This involves data modeling, rich querying capabilities to support real-time decision support and mining, and novel evaluation and optimization processing techniques. In addition, the kind of decision support over data streams is quite different from "traditional" decision-making: decisions are "tactical" rather "strategic". Research on data streams is currently among the most active areas in database research community. We believe that flexible and efficient stream querying will be a crucial component of any future data management and decision support system.

BACKGROUND

The database research community has responded with an abundance of ideas, prototypes and architectures to address the new issues involved in data stream management systems (DSMS). STREAM is Stanford University's approach for a general-purpose DSMS (Arasu et al., 2003); Telegraph and TelegraphCQ (Madden & Franklin, 2002; Chandrasekaran et al., 2003) are prototypes focused on handling measurements of sensor networks, developed in Berkeley; Aurora is a joint project between Brandeis University, Brown University and MIT (Carney et al., 2002) targeted towards stream monitoring applications; AT&T's Hancock (Cortes, Fisher, Pregibon, Rogers & Smith, 2000) and Gigascope (Cranor, Johnson, Spatscheck & Shkapenyuk, 2003) projects are special-purpose data stream systems for network management; Tribeca (Sullivan, 1996) and NiagaraCQ (Chen, DeWitt, Tian & Wang, 2000) are other well-known projects from Telcordia

and University of Wisconsin respectively. The objective of all these projects is to develop systems that can support the challenging analysis requirements of streaming applications.

Furthermore, a plethora of articles, papers and tutorials appeared recently in the research literature. Some of the most well known survey articles follow. Faloutsos (2004), discusses indexing and mining techniques over data streams; Koudas and Srivastava (2003), present the state-of-the-art algorithms on stream query processing; Muthukrishnan (2003), reviews data stream algorithms and applications; Babcock et al. (2002), present an excellent survey on data stream prototypes and issues; Garofalakis, Gehrke and Rastogi (2002), discuss various existing models and mining techniques over data streams.

APPLICATIONS

Stream applications span a wide range of everyday life. Real-time analytics is an essential part of these applications and becomes rapidly more and more critical in decision-making. It is apparent from the list of areas below that efficiently querying and processing data streams is a necessary element of modern decision support systems.

- **Telecommunications.** The telecommunications sector is undoubtedly one of the prime beneficiaries of such data management systems due to the huge amount of data streams that govern voice and data communications over the network infrastructure. Examples of stream analysis include fraud detection, real-time billing, dynamic pricing, network management, traffic monitoring and so on. Streams (calls, packets) have to be mined at real-time to discover outliers and patterns (fraud detection); correlated, joined and aggregated to express complex business rules (billing, dynamic pricing); and monitored – computing averages and min./max. values over periods of time - to uncover unusual traffic patterns (network management).
- **Sensors.** Sensor technology becomes extremely wide-spread and it will probably be the next killer app: large number of cheap, wireless sensors attached to products, cars, computers, even sport players and animals, tracking and digitizing behavior, traffic, location and motion. Examples involve electronic property stickers (super markets, libraries, shopping carts, etc.), vehicle sensors (to pay electronically tolls, route traffic, set speed, etc.) and location-identification sensors (to report location, serve content, detect routes, etc.) A “sensor” world leads to a “stream” world. Millions of input data

every few seconds need to be analyzed: aggregate (what is the average traffic speed), correlate (two products sell together), alert (quantity of a product is below a threshold), localize and monitor.

- **Finance.** Financial data streams come in many different forms: stock tickers, news feeds, trades, etc. Financial companies want to analyze these streams at real-time and take “tactical” business decisions (opposed to “strategic” decisions, associated to OLAP or data mining). For example, Charles Schwab wants to compute commission on each trade at real-time; Fidelity would like to route content in trades at real-time. Traderbot (www.traderbot.com) is a web-based financial search engine that evaluates queries (both traditional and continuous) over real-time streaming data (e.g. “find all stocks between •20 and •200 where the spread between the high tick and the low tick over the past 30 minutes is greater than 3% of the last price and in the last 5 minutes the average volume has surged by more than 300%.”)
- **Web Management.** Large web sites monitor web logs (clickstreams) online to enable applications such as personalization, performance monitoring, and load balancing. Some web sites served by widely distributed web servers (e.g. Yahoo) may need to coordinate many distributed clickstream analyses, e.g. to track heavily accessed web pages (e.g. CNN, BBC) as part of their real-time performance monitoring. (Babcock et al., 2002)
- **Network Management.** Network traffic management systems monitor a variety of continuous data streams at real-time, such as packet traces, packet flows and performance measurements in order to compute statistics, detect anomalies, adjust routing, etc. The volume of data streams can be humongous and thus, query processing must be done very carefully.
- **Military.** One of the most interesting applications in military is battalion monitoring – where sensors are installed on every vehicle, human, etc. – having thousands of sensors reporting state in real-time. In these applications we want to know each time where vehicles and troops are. Examples include queries such as “tell me when 3 of my 4 tanks have crossed the front line” and “tell me when someone is pointing a gun at me”.

ISSUES AND CHALLENGES

Performing decision-making queries on top of data streams is a major challenge. For example, only one-pass algorithms are allowed (because data can be seen only once) and memory has to be managed very carefully (what to keep and what to discard). The need for a data stream

7 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage:

www.igi-global.com/chapter/incorporating-data-stream-analysis-into/14451

Related Content

Information Overload as a Challenge and Changing Point for Educational Media Literacies

Sonja Ganguin, Johannes Gemkow and Rebekka Haubold (2017). *Information and Communication Overload in the Digital Age* (pp. 302-328).

www.irma-international.org/chapter/information-overload-as-a-challenge-and-changing-point-for-educational-media-literacies/176576

The Social Context of Knowledge

Daniel Memmi (2009). *Selected Readings on Information Technology Management: Contemporary Issues* (pp. 21-39).

www.irma-international.org/chapter/social-context-knowledge/28659

Quantifying the Risk of Intellectual Property Loss in Analytics Outsourcing

Handanhal Ravinder, Ram B. Misra and Haiyan Su (2015). *Information Resources Management Journal* (pp. 1-16).

www.irma-international.org/article/quantifying-the-risk-of-intellectual-property-loss-in-analytics-outsourcing/125894

The Influence of Organization Structure and Organizational Learning Factors on the Extent of EDI Implementation in U.S. Firms

Matthew K. McGowan and Gregory R. Madey (1998). *Information Resources Management Journal* (pp. 17-27).

www.irma-international.org/article/influence-organization-structure-organizational-learning/51053

Discrete Total Variation-Based Non-Local Means Filter for Denoising Magnetic Resonance Images

Nikita Joshi, Sarika Jain and Amit Agarwal (2020). *Journal of Information Technology Research* (pp. 14-31).

www.irma-international.org/article/discrete-total-variation-based-non-local-means-filter-for-denoising-magnetic-resonance-images/264755