

Chapter 4

Big Data Time Series Stream Data Segmentation Methods

Dima Alberg

Shamoon College of Engineering (SCE), Israel

ABSTRACT

In this chapter, the authors introduce the interval sliding window (ISW) and interval sliding window and bottom up (ISWAB) algorithms, which are applicable to big data numerical time series data streams and use as input the confidence level parameter rather than the maximum error threshold. The proposed algorithms have two advantages: first, they allow performance comparisons across different time series data streams without changing the algorithm settings, and second, they do not require preprocessing the original time series data stream in order to determine heuristically the reasonable error value. These improvements are very efficient and important in context of big data time series data streams processing. Finally, an empirical evaluation was performed on two types of time series data.

INTRODUCTION

Big data time series data streams are ubiquitous in finance, meteorology and engineering. It may be impossible to process an entire “big data” continuous data stream or to scan through it multiple times due to its tremendous volume. In Heraclitus’s well-known saying, “*you never step in the same stream twice,*” and so it is with “big data” temporal data streams.

Unlike traditional data sets, big data continuous data streams flow into a computer system continuously, in a non-stationary way and with varying update rates. They are time-stamped, fast-changing, massive, and potentially infinite.

Under these circumstances, they represent an application area of growing importance in the data mining research. For example, sensors generate one million samples every minute therefore the primary purpose of time series data stream segmentation is dimensionality reduction. This technique is used in many areas of data stream mining as: frequent patterns finding, structural changes and concept drifts detection (Tabassum & Gama, 2016), time series classification and prediction (Hulten & Domingos, 2003), time series similarities searching (Mori, Mendiburu, Keogh & Lozano, 2016) etc. The main

DOI: 10.4018/978-1-5225-7598-6.ch004

principle of segmentation algorithms concludes in reducing the big data time series dimensionality by dividing the time axis into intervals behaving approximately according to a simple model. A good big data time series data stream segmentation algorithm must be OFASC (Online, Fast, Accurate, Simple and Comparable). For example the Sliding Window algorithm (Keogh, Chu, Hart, & Pazzani, 2004) on the one hand is online (O), very fast (F) and relatively simple (S) for using in online segmentation applications but on the other hand, it sometimes gives poor accuracy (A) and does not allow to perform online multivariate segmentation (C). Therefore, we will classify this algorithm to OFS segmentation algorithms domain.

The segmentation problem can be defined in following way: first, given a time series data stream to produce the best representation such that the maximum error for any segment does not exceed some user specified confidence level error threshold. It is important to add, that using a relative parameter such as confidence level will allow to evaluate an online multivariate segmentation and second, to construct a user friendly segmentation application which will evaluate and compare the proposed online segmentation algorithms in real time. As we shall see in later sections, the state-of-the-art segmentation algorithms do not meet all these requirements.

The rest of the paper is organized as follows. In Section 2, we provide a literature review of three state-of-the-art online piecewise linear segmentation algorithms. In Section 3, we provide a methodology for improving the existing state-of-the-art online segmentation algorithms. The proposed methodology based on novel bound error estimation, which uses a relative probability parameter instead of maximum error nominal parameter and meets the proposed OFASC requirements. Section 4 briefly demonstrates a real-time segmentation application. Finally, in Section 5 and 6 we provide brief and meaningful empirical comparison of the proposed algorithms and suggest final conclusions.

BACKGROUND

Several high level representations of time series have been proposed in the research literature, including Fourier Transforms (Keogh et al., 2000), Wavelets (Chan & Fu, 1999), Symbolic Mappings (Das, Lin, Mannila, Renganathan, & Smyth, 1998; Perng et al., 2000) and Piecewise Linear Approximation or PLA: (Chan & Fu, 1999; Ge & Smyth, 1999; Hunter & McIntosh, 1998; Junker, Amft, Lukowicz, & Tröster, 2008; Keogh et al., 2004; Lavrenko, Schmill, Lawrie, Ogilvie, Jensen, & Allan, 2000; Li, Yu, & Castelli, 1998; Osaki, Shimada, & Uehara, 1999; Park, Lee, & Chu, 1999; Qu, Wang, & Wang, 1998; Shatkay & Zdonik, 1996; Vullings, Verhaegen, & Verbruggen, 1997; Wang & Wang, 2000).

In this work, our attention will confine to PLA, perhaps the most frequently used representation in continuous time series data streams. Obviously, all piecewise linear segmentation algorithms can also be classified as batch or online (Vullings et al., 1997). The problem discussed by (Keogh et al., 2004) is actually how to build online, fast and accurate algorithm for piecewise linear segmentation of time series data stream, because on the one hand, the main problem of online Sliding Window algorithm (Keogh et al., 2004) concerns in its poor accuracy (Qu et al., 1998; Wang & Wang, 2000) and its inability to look ahead. On the other hand the offline accurate Bottom Up (Keogh et al., 2004) algorithm is impractical or may even be unfeasible in a data mining context, where the data are in the order of terabytes or arrive in continuous streams. This problem is very important because for scalability purposes the proposed piecewise linear segmentation algorithm needs to capture the online nature of sliding windows and yet retain the superiority of Bottom Up.

9 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/big-data-time-series-stream-data-segmentation-methods/214603

Related Content

Headache App: Usability Assessment and Criterion Validity

Tânia Dantas, Milton Rodrigues dos Santos, Alexandra Queirós and Anabela G. Silva (2018). *International Journal of Mobile Computing and Multimedia Communications* (pp. 1-11).

www.irma-international.org/article/headache-app/205676

Information Flow Control Based on the CapBAC (Capability-Based Access Control) Model in the IoT

Shigenari Nakamura, Tomoya Enokido and Makoto Takizawa (2019). *International Journal of Mobile Computing and Multimedia Communications* (pp. 13-25).

www.irma-international.org/article/information-flow-control-based-on-the-capbac-capability-based-access-control-model-in-the-iot/241785

What is New about the Internet Delay Space?

Zhang Guomin, Wang Zhanfeng, Wang Rui, Wang Na and Xing Changyou (2014). *International Journal of Mobile Computing and Multimedia Communications* (pp. 36-55).

www.irma-international.org/article/what-is-new-about-the-internet-delay-space/144444

Customer Relationship Management on Internet and Mobile Channels: An Analytical Framework and Research Directions

Susy S. Chan and Jean Lam (2009). *Mobile Computing: Concepts, Methodologies, Tools, and Applications* (pp. 2212-2232).

www.irma-international.org/chapter/customer-relationship-management-internet-mobile/26661

Interaction Design for Personal Photo Management on a Mobile Device

Hyowon Lee, Cathal Gurrin, Gareth J.F. Jones and Alan F. Smeaton (2008). *Handbook of Research on User Interface Design and Evaluation for Mobile Technology* (pp. 69-85).

www.irma-international.org/chapter/interaction-design-personal-photo-management/21824