

Detecting Online Conversations Going Viral

Time-series aware evaluation of change detection
algorithms



Matt Chapman

matthew.chapman@student.uva.nl

Spring 2017, 12 pages

Supervisor: Evangelos Kanoulas, Universiteit van Amsterdam
Host organisation: Buzzcapture International, <http://www.buzzcapture.com>



UNIVERSITEIT VAN AMSTERDAM
FACULTEIT DER NATUURWETENSCHAPPEN, WISKUNDE EN INFORMATICA
MASTER SOFTWARE ENGINEERING
<http://www.software-engineering-amsterdam.nl>

Contents

Abstract	2
1 Problem Statement & Motivation	3
1.1 Problem Statement	3
1.2 Motivation	4
2 Research Method	5
2.1 Data Preparation	5
2.2 Scoring Metrics	5
3 Background & Context	6
4 Research	7
5 Results	8
6 Analysis & Conclusions	9

Abstract

Write abstract

Chapter 1

Problem Statement & Motivation

1.1 Problem Statement

Within the domain of change detection algorithms there are a number of available methods for evaluating the efficiency and accuracy of a given algorithm, depending on how the problem is framed. The different problem framings available are, for example:

Classification Wherein the algorithm result fits into one of two (or more) classes, for example correct or incorrect.

are these definitions correct?

Clustering Wherein the algorithm results are formed into clusters and evaluated using clustering metrics.

Partitioning Like clustering, but?

Retrieval Wherein the results of the algorithm are scored as a *retrieval problem*, where detection relevance is taken into account, perhaps along with some form of temporal or redundancy penalty.

How is this different from clustering?

With this research, it is intended to address this dichotomy between approaches, and further suggest an evaluation approach that is suitable for use when attempting to evaluate change detection algorithms applied to time-series data specifically.

While change detection itself has been around since the 1930's (and *online* change detection since the 1950's) it is still a field that attracts new thoughts and approaches - one example of which being Ginsberg et al.'s "Detecting influenza epidemics using search engine query data." [Gin+09].

Various attempts have been made to evaluate the myriad approaches to change detection ([BNZ14] for example) but each of these attempts tend to frame the problem somewhat differently, or do not take into account change detection in time-series data.

This research intends to address the following research questions:

RQ1 Are there deficiencies in existing methods for evaluating change detection algorithms?

RQ2 Do certain measures perform better as an evaluation method when applied against changes detected in a data stream with certain properties?

RQ3 Is there room for adjustment in existing measures, such that they can be made more effective for the evaluation of change detection algorithms?

this needs expansion, I think

1.2 Motivation

Change detection first came about as a quality control measure in manufacturing, and methods within this domain are generally referred to as *control charts*. Since the inception of approaches such as CUSUM that provide the possibility for on-line evaluation of continuous data streams, change detection has grown as a field. With applications such as epidemic detection, online reputation management and infrastructure error detection, change detection is hugely useful.

This particular research is motivated specifically by the online reputation management sector. The business hosting this research project (Buzzcapture International [<http://www.buzzcapture.com>]) is a Dutch online reputation management company that provides services to other businesses throughout europe. Chief among these is the BrandMonitor application, which, among other features, provides a rudimentary notification system for clients that is triggered once there is an increase in conversation volume of $\%n$. It is the intention of this research to provide a robust evaluation method for change detection algorithms such that an approach that is most effective for this particular use case can be selected and implemented.

Chapter 2

Research Method

2.1 Data Preparation

2.2 Scoring Metrics

Example function:

$$f(\textit{relevance}, \textit{temporal_penalty}, \textit{redundancy_penalty}) \quad (2.1)$$

start designing scoring methods

Possible relevance measure, where t_0 is the earliest a spike can be detected and t_n is the time that the signal returns to normal. $f(x)$ describes the function of the curve:

$$\int_{t_n}^{t_0} f(x) dx \quad (2.2)$$

start designing relevance measure

planning to do something with the area under the curve for relevance

Chapter 3

Background & Context

This research was primarily borne out of reading “A Brief Comparison of Algorithms for Detecting Change Points in Data” by Buntain, Natoli, and Zivkovic[\[BNZ14\]](#)

Chapter 4

Research

Chapter 5

Results

Chapter 6

Analysis & Conclusions

Bibliography

- [AF13] Leman Akoglu and Christos Faloutsos. “Anomaly, event, and fraud detection in large network datasets”. In: *Proceedings of the sixth ACM international conference on Web search and data mining - WSDM '13* (2013), p. 773. DOI: [10.1145/2433396.2433496](https://doi.org/10.1145/2433396.2433496). URL: <http://dl.acm.org/citation.cfm?id=2433396.2433496>.
- [Alv+11] Foteini Alvanaki et al. “EnBlogue: emergent topic detection in Web 2.0 streams”. In: *Proc. ACM SIGMOD International Conference on Management of Data* (2011), pp. 1271–1274. ISSN: 07308078. DOI: [10.1145/1989323.1989473](https://doi.org/10.1145/1989323.1989473). URL: <http://doi.acm.org/10.1145/1989323.1989473%7B%5C%7D5Cnpapers2://publication/uuid/44E009D1-8574-49C1-A0AC-C2B247FE9043>.
- [Ami+09] Enrique Amigó et al. “A comparison of extrinsic clustering evaluation metrics based on formal constraints”. In: *Information Retrieval* 12.4 (2009), pp. 461–486. ISSN: 13864564. DOI: [10.1007/s10791-008-9066-8](https://doi.org/10.1007/s10791-008-9066-8).
- [BN93] M Basseville and Igor V Nikiforov. *Detection of Abrupt Changes: Theory and Application*. 1993. ISBN: 0-13-126780-9. DOI: [10.1016/0967-0661\(94\)90196-1](https://doi.org/10.1016/0967-0661(94)90196-1). URL: <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.77.6896%7B%5C%7Drep1%7B%5C%7Dtype=pdf>.
- [BPP07] S. Bersimis, S. Psarakis, and J. Panaretos. “Multivariate statistical process control charts: An overview”. In: *Quality and Reliability Engineering International* 23.5 (2007), pp. 517–543. ISSN: 07488017. DOI: [10.1002/qre.829](https://doi.org/10.1002/qre.829).
- [BNZ14] Cody Buntain, Christopher Natoli, and Miroslav Zivkovic. “A Brief Comparison of Algorithms for Detecting Change Points in Data”. In: *Supercomputing*. 2014. URL: <https://github.com/cbuntain/ChangePointDetection>.
- [Das+09] Tamraparni Dasu et al. “Change (detection) you can believe in: Finding distributional shifts in data streams”. In: *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)* 5772 LCNS (2009), pp. 21–34. ISSN: 03029743. DOI: [10.1007/978-3-642-03915-7_3](https://doi.org/10.1007/978-3-642-03915-7_3).
- [DDD05] Frédéric Desobry, Manuel Davy, and Christian Doncarli. “An online Kernel change detection algorithm”. In: *IEEE Transactions on Signal Processing* 53.8 (2005), pp. 2961–2974. ISSN: 1053587X. DOI: [10.1109/TSP.2005.851098](https://doi.org/10.1109/TSP.2005.851098).
- [Dow08] Allen B. Downey. “A novel changepoint detection algorithm”. In: *Applied Microbiology and Biotechnology* (2008), pp. 1–11. arXiv: [0812.1237](https://arxiv.org/abs/0812.1237). URL: <http://arxiv.org/abs/0812.1237>.
- [FP99] Tom Fawcett and Foster Provost. “Activity monitoring: Noticing interesting changes in behavior”. In: *Proceedings of the fifth ACM SIGKDD international conference on Knowledge discovery and data mining* 1.212 (1999), pp. 53–62. ISSN: 09258574. DOI: [10.1016/j.ecoleng.2010.11.031](https://doi.org/10.1016/j.ecoleng.2010.11.031). URL: <http://portal.acm.org/citation.cfm?id=312195>.
- [GP04] Pedro Galeano and Daniel Peña. “Variance Changes Detection in Multivariate Time Series”. 2004.
- [Gin+09] Jeremy Ginsberg et al. “Detecting influenza epidemics using search engine query data.” In: *Nature* 457.7232 (2009), pp. 1012–4. ISSN: 1476-4687. DOI: [10.1038/nature07634](https://doi.org/10.1038/nature07634). URL: <http://www.ncbi.nlm.nih.gov/pubmed/19020500>.

- [KS09] Yoshinobu Kawahara and Masashi Sugiyama. “Change-point detection in time-series data by direct density-ratio estimation”. In: *Proceedings of the 2009 SIAM International Conference on Data Mining* (2009), pp. 389–400.
- [KBG04] Daniel Kifer, Shai Ben-david, and Johannes Gehrke. “Detecting Change in Data Streams”. In: *Proceedings of the 30th VLDB Conference* (2004), pp. 180–191. ISSN: 00394564. DOI: [10.1016/0378-4371\(94\)90421-9](https://doi.org/10.1016/0378-4371(94)90421-9). arXiv: [9310008](https://arxiv.org/abs/9310008) [cond-mat].
- [Kul+05] Martin Kulldorff et al. “A space-time permutation scan statistic for disease outbreak detection”. In: *PLoS Medicine* 2.3 (2005), pp. 0216–0224. ISSN: 15491277. DOI: [10.1371/journal.pmed.0020059](https://doi.org/10.1371/journal.pmed.0020059).
- [LS99] Tze Leung Lai and Jerry Zhaolin Shan. “Efficient recursive algorithms for detection of abrupt changes in signals and control systems”. In: *IEEE Transactions on Automatic Control* 44.5 (1999), pp. 952–966. ISSN: 00189286. DOI: [10.1109/9.763211](https://doi.org/10.1109/9.763211).
- [MJ12] David S. Matteson and Nicholas A. James. “A nonparametric approach for multiple change point analysis of multivariate data”. In: *Submitted* 14853 (2012), pp. 1–29. ISSN: 0162-1459. DOI: [10.1080/01621459.2013.849605](https://doi.org/10.1080/01621459.2013.849605). arXiv: [1306.4933](https://arxiv.org/abs/1306.4933).
- [PRG10] Anita M Pelecanos, Peter a Ryan, and Michelle L Gatton. “Outbreak detection algorithms for seasonal disease data: a case study using Ross River virus disease.” In: *BMC medical informatics and decision making* 10.1 (2010), p. 74. ISSN: 1472-6947. DOI: [10.1186/1472-6947-10-74](https://doi.org/10.1186/1472-6947-10-74). URL: <http://www.biomedcentral.com/1472-6947/10/74>.
- [Qah+15] Abdulhakim A. Qahtan et al. “A PCA-Based Change Detection Framework for Multidimensional Data Streams”. In: *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining - KDD '15* (2015), pp. 935–944. DOI: [10.1145/2783258.2783359](https://doi.org/10.1145/2783258.2783359). URL: <http://dl.acm.org/citation.cfm?id=2783258.2783359>.
- [SV95] D. Siegmund and E. S. Venkatraman. “Using the generalized likelihood ratio statistic for sequential detection of a change-point”. In: *The Annals of Statistics* 23.1 (1995), pp. 255–271. ISSN: 0090-5364. DOI: [10.1214/aos/1176324466](https://doi.org/10.1214/aos/1176324466).
- [TR05] Alexander G Tartakovsky and Boris L Rozovskii. “A Nonparametric Multichart CUSUM Test for Rapid Intrusion Detection”. In: *Proceedings of Joint Statistical Meetings* (2005), pp. 7–11.
- [Tar+06] Alexander G. Tartakovsky et al. “Detection of intrusions in information systems by sequential change-point methods”. In: *Statistical Methodology* 3.3 (2006), pp. 252–293. ISSN: 15723127. DOI: [10.1016/j.stamet.2005.05.003](https://doi.org/10.1016/j.stamet.2005.05.003).
- [TGS14] Dang-Hoan Tran, Mohamed Medhat Gaber, and Kai-Uwe Sattler. “Change Detection in Streaming Data in the Era of Big Data: Models and Issues”. In: *ACM SIGKDD Explorations Newsletter - Special issue on big data* 1 (2014), pp. 30–38. ISSN: 1931-0145. DOI: [10.1145/2674026.2674031](https://doi.org/10.1145/2674026.2674031).
- [WJ76] Alan S. Willsky and Harold L. Jones. “A Generalized Likelihood Ratio Approach to the Detection and Estimation of Jumps in Linear Systems”. In: *IEEE Transactions on Automatic Control* 21.1 (1976), pp. 108–112. ISSN: 15582523. DOI: [10.1109/TAC.1976.1101146](https://doi.org/10.1109/TAC.1976.1101146).
- [Xu+11] Yi Xu et al. “Pattern change discovery between high dimensional data sets”. In: *Proceedings of the 20th ACM international conference on Information and knowledge management - CIKM '11* (2011), p. 1097. DOI: [10.1145/2063576.2063735](https://doi.org/10.1145/2063576.2063735). URL: <http://dl.acm.org/citation.cfm?doid=2063576.2063735>.

ToDo Notes

Write abstract	2
are these definitions correct?	3
How is this different from clustering?	3
this needs expansion, I think	3
start designing scoring methods	5
planning to do something with the area under the curve for relevance	5
start designing relevance measure	5