

CP based Sequence Mining on the cloud using spark

Dissertation presented by
Cyril DE VOGELAERE

for obtaining the Master's degree in
Computer Science

Supervisor(s)
Pierre SCHAUS

Reader(s)
John AOGA, Guillaume Derval

Academic year 2016-2017

Contents

1	Introduction	4
2	Sequential Pattern Mining	4
2.1	Sequential Pattern Mining Background	4
2.1.1	Definitions and Concepts	4
2.1.2	Sequences of Symbols VS Sequences of Set of Symbols	4
2.2	Existing specialised approaches	4
2.2.1	Prefix-Span	4
2.2.2	apriori	4
2.2.3	GSP	4
2.2.4	cspade	4
2.2.5	FP-Growth	4
2.3	Existing CP Based approaches	5
2.3.1	CPSM	5
2.3.2	PP	5
2.3.3	Gap-Seq	5
2.3.4	PPIC	5
2.4	Parallelisation	5
2.4.1	The Benefits of Parallelisation	5
2.4.2	Tool selection	5
3	Implementation of a Scalable CP Based Algorithm	5
3.1	Spark's original implementation	5
3.2	A First Scalable CP based Implementation	5
3.3	Option for further improvements	5
3.4	Improving the Switch to a CP Local Execution	5
3.4.1	Quicker - Start	5
3.4.2	Cleaning Sequence before the Local Execution	5
3.5	Improving the Scalable Execution	5
3.5.1	Automatic Choice for Local Execution	5
3.5.2	Position lists	5
3.5.3	Specialising the Scalable Execution	5
3.5.4	Priority Scheduling for Sub-Problems	5
3.6	CP Based Local Execution for Sequence of Sets of Symbols	5
3.6.1	Pushing PPIC's Ideas Further	5
3.6.2	Adding Partial Projections to PPIC	5
4	Performances	5
4.1	Datasets & Number of Partitions	5
4.2	Performance Testing Procedure	5
4.2.1	Distribution Choice & Cluster Architecture	5
4.2.2	Program Parameters	5
4.2.3	Measurement Span	5
4.3	Testing the Implementations	5
4.4	Scalability Tests	5
5	Conclusion	5
6	Annexes	5
	References	5
	Acronyms	7

List of Figures

List of Tables

List of Algorithms

Abstract

TODO - Half a page to a page of content should be enough

1 Introduction

TODO Prefix Projection Incremental Counting propagator mathematics Sequence of Symbols
Sequence of Sets of Symbols

2 Sequential Pattern Mining

2.1 Sequential Pattern Mining Background

2.1.1 Definitions and Concepts

2.1.2 Sequences of Symbols VS Sequences of Set of Symbols

2.2 Existing specialised approaches

2.2.1 Prefix-Span

2.2.2 apriori

2.2.3 GSP

2.2.4 cspade

2.2.5 FP-Growth

Should I talk about FP-Growth ? you didn't mention it in your sheet of paper

2.3 Existing CP Based approaches

2.3.1 CPSM

2.3.2 PP

2.3.3 Gap-Seq

2.3.4 PPIC

2.4 Parallelisation

2.4.1 The Benefits of Parallelisation

2.4.2 Tool selection

3 Implementation of a Scalable CP Based Algorithm

3.1 Spark's original implementation

3.2 A First Scalable CP based Implementation

3.3 Option for further improvements

3.4 Improving the Switch to a CP Local Execution

3.4.1 Quicker - Start

3.4.2 Cleaning Sequence before the Local Execution

3.5 Improving the Scalable Execution

3.5.1 Automatic Choice for Local Execution

3.5.2 Position lists

3.5.3 Specialising the Scalable Execution

3.5.4 Priority Scheduling for Sub-Problems

3.6 CP Based Local Execution for Sequence of Sets of Symbols

3.6.1 Pushing PPIC's Ideas Further

3.6.2 Adding Partial Projections to PPIC

4 Performances

4.1 Datasets & Number of Partitions

4.2 Performance Testing Procedure

4.2.1 Distribution Choice & Cluster Architecture

4.2.2 Program Parameters

4.2.3 Measurement Span

4.3 Testing the Implementations

4.4 Scalability Tests

5 Conclusion

6 Annexes

References

- [1] M. Zaharia, M. Chowdhury, T. Das, A. Dave, J. Ma, M. McCauley, M. J. Franklin, S. Shenker, and I. Stoica, “Resilient distributed datasets: A fault-tolerant abstraction for in-memory cluster computing,” in *Proceedings of the 9th USENIX conference on Networked Systems Design and Implementation*, pp. 2–2, USENIX Association, 2012.
- [2] Oscala Team, “Oscala: Scala in OR,” 2012. Available from <https://bitbucket.org/oscarlib/oscar>.
- [3] N. R. Mabroukeh and C. I. Ezeife, “A taxonomy of sequential pattern mining algorithms,” *ACM Computing Surveys (CSUR)*, vol. 43, no. 1, p. 3, 2010.
- [4] A. Kemmar, S. Loudni, Y. Lebbah, P. Boizumault, and T. Charnois, “Prefix-projection global constraint for sequential pattern mining,” in *International Conference on Principles and Practice of Constraint Programming*, pp. 226–243, Springer, 2015.
- [5] A. Kemmar, S. Loudni, Y. Lebbah, P. Boizumault, and T. Charnois, “A global constraint for mining sequential patterns with gap constraint,” in *International Conference on AI and OR Techniques in Constraint Programming for Combinatorial Optimization Problems*, pp. 198–215, Springer, 2016.
- [6] T. Guns, S. Nijssen, and L. De Raedt, “Itemset mining: A constraint programming perspective,” *Artificial Intelligence*, vol. 175, no. 12–13, pp. 1951–1983, 2011.
- [7] J. O. Aoga, T. Guns, and P. Schaus, “An efficient algorithm for mining frequent sequence with constraint programming,” in *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pp. 315–330, Springer, 2016.
- [8] J. Pei, J. Han, B. Mortazavi-Asl, J. Wang, H. Pinto, Q. Chen, U. Dayal, and M.-C. Hsu, “Mining sequential patterns by pattern-growth: The prefixspan approach,” *IEEE Transactions on knowledge and data engineering*, vol. 16, no. 11, pp. 1424–1440, 2004.
- [9] Z. Yang, Y. Wang, and M. Kitsuregawa, “Lapin: effective sequential pattern mining algorithms by last position induction for dense databases,” *Advances in Databases: Concepts, Systems and Applications*, pp. 1020–1023, 2007.
- [10] Z. Yang and M. Kitsuregawa, “Lapin-spam: An improved algorithm for mining sequential pattern,” in *Data Engineering Workshops, 2005. 21st International Conference on*, pp. 1222–1222, IEEE, 2005.
- [11] F. Bonchi and C. Lucchese, “Extending the state-of-the-art of constraint-based pattern discovery,” *Data & Knowledge Engineering*, vol. 60, no. 2, pp. 377–399, 2007.
- [12] L. De Raedt, T. Guns, and S. Nijssen, “Constraint programming for itemset mining,” in *Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 204–212, ACM, 2008.
- [13] L. De Raedt, T. Guns, and S. Nijssen, “Constraint programming for data mining and machine learning,” in *Proceedings of the Twenty-Fourth AAAI Conference on Artificial Intelligence (AAAI-10)*, pp. 1671–1675, 2010.
- [14] L. D. Raedt and A. Zimmermann, “Constraint-based pattern set mining,” in *proceedings of the 2007 SIAM International conference on Data Mining*, pp. 237–248, SIAM, 2007.
- [15] M. J. Zaki, “Sequence mining in categorical domains: incorporating constraints,” in *Proceedings of the ninth international conference on Information and knowledge management*, pp. 422–429, ACM, 2000.

Acronyms

PPIC Prefix Projection Incremental Counting propagator. 4

SoS Sequence of Symbols. 4

SoSS Sequence of Sets of Symbols. 4

