

Auto-Scaling Techniques for Spark Streaming

Master-Thesis von Seyedmajid Azimi Gehraz

Tag der Einreichung:

1. Gutachten: Prof. Dr. rer. nat. Carsten Binnig
2. Gutachten: Dr. Thomas Heinze



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Fachbereich Informatik
Data Management

Auto-Scaling Techniques for Spark Streaming

Vorgelegte Master-Thesis von Seyedmajid Azimi Gehraz

1. Gutachten: Prof. Dr. rer. nat. Carsten Binnig
2. Gutachten: Dr. Thomas Heinze

Tag der Einreichung:

Contents

List of Figures	II
List of Tables	III
1 Abstract	1
2 Introduction to Auto-Scaling	2
3 Spark Streaming	3
4 Reinforcement Learning	4
5 Design	5
6 Implementation Detail	6
7 Evaluation	7
8 Related Work	8
9 Conclusion	10
Bibliography	IV

List of Figures

List of Tables



3 Spark Streaming



4 Reinforcement Learning

6 Implementation Detail

Dynamic resource scaling in cloud environments has been studied extensively in literature. As a naive implementation, there are two thresholds, namely *upper bound* and *lower bound*. However, such an implementation suffers from *oscillating* decisions. In order to remedy this issue, *grace period* shall be enforced. During this period, no scaling decision is made.

Hasan et al. [4] introduced four thresholds and two time periods. *ThrUpper* defines upper bound. *ThrBelowUpper* is slightly below *ThrUpper*. Similarly, *ThrLower* defines lower bound and *ThrAboveLower* is slightly above the lower bound. In case, system utilization stays between *ThrUpper* and *ThrBelowUpper* for a specific duration, then cluster controller decides to take a scale-out action, by adding resources. On the other hand, if system utilization stays between *ThrLower* and *ThrAboveLower* for a specified duration, then controller decides to take scale-in action. Defining two levels of thresholds helps to detect workload *persistence* and avoid making immature scaling decision. However, defining thresholds is a tricky and manual process, and need to be carefully done[1]. It should be noted that, computation overhead of this approach is very low.

RightScale applies voting algorithm[12] among nodes to make scaling decision. In order for a specific action to be decided, majority of nodes should vote in favour of that action. Otherwise, no action is selected as a default action. Afterwards, nodes apply grace period to stabilize the cluster. The complexity of the voting process in trusted environments is in the order of $O(n^2)$, which leads to heavy network traffic among participants when cluster size grows. This approach is also categorized in threshold based approaches. Thus, it suffers from the same issue as mentioned above.

Herbst et al. [7] surveys different auto-scaling techniques based on time-series analysis in order to forecast *trends* and *seasons*. *Moving Average Method* takes the average over a sliding window and smooths out minor noise level. Its computational overhead is proportional to size of the window. *Simple Exponential Smoothing*(SES) goes further than just taking average. It gives more weight to more recent values in sliding window by an exponential factor. Although it is more computationally intensive compared to moving average, it is still negligible. SES is capable of detecting short-term trends but fails at predicting seasons. These approaches are more specific instances of *ARIMA* (Auto-Regressive Integrated Moving Average) which is a general purpose framework to calculate moving averages. However, time-series analysis is only suitable for stationary problems consist of recurring workload patterns such as web applications. In case of streaming (specially in multi-tenant environments), in which the workload is highly unpredictable, time-series analysis tends to be less useful. Additionally, more advanced forms of time-series analysis which are capable of forecasting seasons (such as *tBATS Innovation State Space Modelling Framework*[10], *ARIMA Stochastic Process Modelling Framework*[9]) are computationally infeasible for streaming workloads.

Herbst et al. [6] continues the survey on state of the art techniques to predict future workload. It includes workload forecasting based on *Bayesian Networks*(BN) and *Neural Networks*. There are issue with each of them that makes them unsuitable for streaming workloads. As an example, there is no universally applicable method to construct a BN. Furthermore, it requires collecting data and training the model offline. Neural networks suffer from the same issue. That is, it requires collecting samples and training the model offline. For complex models, training phase is typically computationally infeasible which was one the requirements.

Tesauro et al. [14] proposes a hybrid approach to overcome poor performance of online training. The system consists of two components: an online component based on queuing system combined with reinforcement learning component that is trained offline. The offline component is based on *neural networks*. Author models the data center as multiple applications managed under a single resource manager. Modelling streaming workloads as a queuing system has two problems. First, modelling is a complicated process and determining probability distributions requires domain knowledge. Second, it requires access to each node (so it can be modelled as a queue) which is currently not possible without modifying spark-core package. Since, it was one the requirements to provide a solution without making any modification to spark-core, this work has been abandoned.

Rao et al. [11] proposed to use reinforcement learning to manage resources consumed by virtual machines. It employs standard model-free learning, which is known as *temporal difference*[13] or *sarsa* algorithm. The state space consists of metrics collected from virtual machines (CPU, RAM, Network IO, ...). There is no global controller and each node decides based on its own Q-Table. As mentioned in literature, standard temporal difference has a slow convergence speed. In order to speedup bootstrap phase, Q-Table is initialized by values that were obtained during separate supervised training. Since this approach also relies on offline training, it wasn't adopted by this thesis.

Heinze et al. [5] implemented reinforcement learning in the context of FUGU[3] (which is a *data stream processing* framework). Each node, has its own Q-Table and imposes local policy without coordinating to other nodes. This architecture can not be applied in context of spark streaming, since spark abstracts away individual nodes from perspective of applicaiton developer. In order to decrease state space, the author applies two techniques. First, only system utilization is considered. Second, system utilization is discretized using coarse grained steps. To remedy slow convergence, controller enforces *monotonicity constraint*[8]. That is, if controller decides to take scale-out action for a specific utilization, it may not decide scale-in for even worse system utilization. This feature has been adopted by this thesis.

Enda, Enda, and Jim [2] proposed a parallel architecture to reinforcement learning. Standard model-free learning (temporal difference) is used. No global controller is involved and each node decides locally. In order to speed up learning, all nodes, maintain two Q-Tables (local and global tables). Local table is learned and updated by each node. Whenever, an agent learns a new value for a specific state, it broadcasts it to other agents. Global table contains values received from other agents. Additionally, agent prioritize local and global tables by assigning weights to each table. Weights are factors that are defined by application developers. Final decision is the outcome of combining local and global tables. Although each node learns some part of the state space (which may overlap with other nodes), it is not applicable in the context of spark streaming. The assumption in this architecture is that, each node is operating autonomously without intervention of other nodes (such as web servers). In contrast, spark is a centrally managed system. That is, all nodes running spark jobs are supervised by a single master node (probably with couple of backup masters).

9 Conclusion

Bibliography

- [1] X. Dutreilh, A. Moreau, J. Malenfant, N. Rivierre, and I. Truck. “From Data Center Resource Allocation to Control Theory and Back”. In: *2010 IEEE 3rd International Conference on Cloud Computing*. 2010, pp. 410–417. DOI: 10.1109/CLOUD.2010.55.
- [2] B. Enda, H. Enda, and D. Jim. “Applying reinforcement learning towards automating resource allocation and application scalability in the cloud”. In: *Concurrency and Computation: Practice and Experience* 25.12 (2012), pp. 1656–1674. DOI: 10.1002/cpe.2864. URL: <https://onlinelibrary.wiley.com/doi/abs/10.1002/cpe.2864>.
- [3] R. Grandl, G. Ananthanarayanan, S. Kandula, S. Rao, and A. Akella. “Multi-resource Packing for Cluster Schedulers”. In: *SIGCOMM Comput. Commun. Rev.* 44.4 (2014), pp. 455–466. ISSN: 0146-4833. DOI: 10.1145/2740070.2626334. URL: <http://doi.acm.org/10.1145/2740070.2626334>.
- [4] M. Z. Hasan, E. Magana, A. Clemm, L. Tucker, and S. L. D. Gudreddi. “Integrated and autonomic cloud resource scaling”. In: *2012 IEEE Network Operations and Management Symposium* (2012), pp. 1327–1334.
- [5] T. Heinze, V. Pappalardo, Z. Jerzak, and C. Fetzter. “Auto-scaling techniques for elastic data stream processing”. In: *2014 IEEE 30th International Conference on Data Engineering Workshops*. 2014, pp. 296–302. DOI: 10.1109/ICDEW.2014.6818344.
- [6] N. Herbst, A. Amin, A. Andrzejak, L. Grunske, S. Kounev, O. J. Mengshoel, and P. Sundararajan. “Online Workload Forecasting”. In: *Self-Aware Computing Systems*. Ed. by S. Kounev, J. O. Kephart, A. Milenkoski, and X. Zhu. Cham: Springer International Publishing, 2017, pp. 529–553. ISBN: 978-3-319-47474-8. DOI: 10.1007/978-3-319-47474-8_18. URL: https://doi.org/10.1007/978-3-319-47474-8_18.
- [7] N. R. Herbst, N. Huber, S. Kounev, and E. Amrehn. “Self-adaptive Workload Classification and Forecasting for Proactive Resource Provisioning”. In: *Proceedings of the 4th ACM/SPEC International Conference on Performance Engineering*. ICPE ’13. Prague, Czech Republic: ACM, 2013, pp. 187–198. ISBN: 978-1-4503-1636-1. DOI: 10.1145/2479871.2479899. URL: <http://doi.acm.org/10.1145/2479871.2479899>.
- [8] H. Herodotou and S. Babu. “Profiling, What-if Analysis, and Cost-based Optimization of MapReduce Programs”. In: *Proceedings of the VLDB Endowment*. Vol. 4. Jan. 2011, pp. 1111–1122.
- [9] R. Hyndman and Y. Khandakar. “Automatic Time Series Forecasting: The forecast Package for R”. In: *Journal of Statistical Software, Articles* 27.3 (2008), pp. 1–22. ISSN: 1548-7660. DOI: 10.18637/jss.v027.i03. URL: <https://www.jstatsoft.org/v027/i03>.
- [10] A. M. D. Livera, R. J. Hyndman, and R. D. Snyder. “Forecasting Time Series With Complex Seasonal Patterns Using Exponential Smoothing”. In: *Journal of the American Statistical Association* 106.496 (2011), pp. 1513–1527. DOI: 10.1198/jasa.2011.tm09771. URL: <https://doi.org/10.1198/jasa.2011.tm09771>.
- [11] J. Rao, X. Bu, C.-Z. Xu, L. Wang, and G. Yin. “VCONF: A Reinforcement Learning Approach to Virtual Machines Auto-configuration”. In: *Proceedings of the 6th International Conference on Autonomic Computing*. ICAC ’09. Barcelona, Spain: ACM, 2009, pp. 137–146. ISBN: 978-1-60558-564-2. DOI: 10.1145/1555228.1555263. URL: <http://doi.acm.org/10.1145/1555228.1555263>.
- [12] RightScale. *Set up Autoscaling using Voting Tags*. Accessed June 20, 2018. 2018. URL: http://support.rightscale.com/12-Guides/Dashboard_Users_Guide/Manage/Arrays/Actions/Set_up_Autoscaling_using_Voting_Tags/.
- [13] R. S. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction*. 2nd ed. The MIT Press, 1998. ISBN: 0262193981.

-
- [14] G. Tesauro, N. K. Jong, R. Das, and M. N. Bennani. “A Hybrid Reinforcement Learning Approach to Autonomic Resource Allocation”. In: *2006 IEEE International Conference on Autonomic Computing*. 2006, pp. 65–73. DOI: 10.1109/ICAC.2006.1662383.