

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
8.1 Introduction	8
8.2 Threshold-Based Techniques	8
8.3 Time-Series Analysis Techniques	8
8.4 Queuing Theory Techniques	9
8.5 Reinforcement Learning Techniques	9
8.6 Summary	11
9 Conclusion	12
Bibliography	IV

List of Figures

List of Tables



3 Spark Streaming



4 Reinforcement Learning

6 Implementation Detail

8 Related Work

8.1 Introduction

Dynamic resource allocation in cloud environments has been studied extensively in literature. In this chapter prior work will be discussed and explored. It is organized as follows. Section 8.2 delves into threshold-based techniques. Section 8.3 investigates techniques based on time-series analysis. Section 8.4 analyses techniques based on queuing theory. Section 8.5 explores Reinforcement Learning techniques comprehensively. Finally, section 8.6 concludes.

8.2 Threshold-Based Techniques

Hasan et al. [10] proposed 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 the central controller decides to take scale-in action. Furthermore, in order to prevent making *oscillating* decisions, *grace period* is enforced. During this period, no scaling decision is made. Defining two levels of thresholds helps to detect workload *persistence* and avoids making immature scaling decision. However, defining thresholds is a tricky and manual process, and needs to be carefully done [6]. It shall be noted that, computation overhead of this approach is very low.

RightScale [25] applies voting algorithm among nodes to make scaling decisions. In order for a specific action to be decided, majority of nodes should vote in favor of that specific action. Otherwise, no-action is elected 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 also suffers from the same issue – accurately adjusting threshold values – as other threshold-based approaches.

Heinze et al. [12] proposed a novel threshold-based solution in the context of FUGU [9] – a data stream processing framework. This techniques uses an adaptive window [1] to monitor the recent changes in workload pattern. In case a change in workload is detected, optimization component is activated and fed with recent short-term utilization history. Thereafter, the optimization component determines monetary cost of current system configuration and then simulates the cost of different scaling decisions. The *latency-aware* cost function has the responsibility to calculate monetary cost of system configuration. The search function is an implementation of *Recursive Random Search* [31] algorithm which consists of two phases. First, in *exploration* phase, the complete parameter space is explored to find a solution with minimum cost. In second phase – *exploitation phase* – only specific parts of the parameter space which has been discovered in first phase, will be investigated. Kielbowicz [17] has implemented this techniques in the context of Spark Streaming. Thus, It is considered in evaluation scenarios.

8.3 Time-Series Analysis Techniques

Herbst et al. [14] 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. Additionally, more advanced forms of time-series analysis which are capable of forecasting seasons (such as *tBATS Innovation State Space Modeling Framework* [19], *ARIMA Stochastic Process Modeling Framework* [16]) are computationally infeasible for streaming workloads.

Taft et al. [28] applied time-series analysis in the context of OLTP databases. The authors argue that reactive approaches don't fit to database world. By the time, auto-scaler component decides to scale-out, it is already too late for a database system. This premise comes from the fact that taking scaling actions in a database doesn't take place in timely manner. The database system has to replicate some of the records which is an additional burden on a heavily loaded system. Thus, database system must take proactive approach and take scaling decisions ahead of time. While this is convincing argument, the auto-scaler module depends on a couple of parameters that are hard to calculate in heterogeneous public cloud environments. First, target throughput of a single server. Second, shortest time to move all database records with single sender-receiver thread. While this might be feasible in some scenarios, on today's cloud environments with virtual machines hosted on heterogeneous physical nodes, getting a near-precise number is unconvincing. It worth noting that author assumed an approximately uniform workload distribution for all database nodes – each database shard serves a fairly equal portion of total workload which is a questionable assumption.

8.4 Queuing Theory Techniques

Lohrmann, Janacik, and Kao [20] proposed a solution based on queuing theory. The solution is designed for *Nephele* [21] streaming engine which has a master-worker style architecture. Similar to Spark Streaming, a job is modeled as a DAG. It utilizes *adaptive output batching* [30] – which is essentially a buffer with variable size – to buffer outgoing messages emitted from one stage to the other. Each task – an executor that runs user defined function (UDF) – is modeled as a G/G/1 queue. That is, the probability distributions of message inter-arrival and service time are unknown. In order to approximate these distributions, a formula proposed by Kingman [18] is used. From a bird's eye view, this solution seems promising. However, authors made two inconceivable assumptions that led us to abandon the proposal. First, worker nodes shall be homogeneous in terms of processing power and network bandwidth. Second, there should be an effective partitioning strategy in place in order to load balance outgoing messages between stages. In reality both assumptions rarely occur. Large scale stream processing clusters are built incrementally. Depending on workload, data skew does exist and imperfect hash functions are widely used by software developers.

Zhang, Cherkasova, and Smirni [32] proposed a solution for multi-tiered enterprise applications based on regression techniques. Regression based models can absorb some level of uncertainty and noise by compacting samples. Each tier is modeled as G/G/1 queue and scaled differently compared to other tiers. The system has fixed number of users – a principle known as *closed-loop queuing network*. In order to calculate system workload – incoming message rate – and service time which is required by queuing models, the authors proposed to use Mean Value Analysis [23]. In order to simplify the queuing network, the system is modeled as a *transaction-based* system with independent requests coming from clients. However, It is widely believed that multi-tiered enterprise applications are *session-based* systems [5]. Each request from the same client depends on her previous request during a specific session.

8.5 Reinforcement Learning Techniques

Herbst et al. [13] surveys on state of the art techniques to predict future workload. It includes workload forecasting based on *Bayesian Networks* (BN) and *Neural Networks*. There are several issues 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 issues. That is, it requires collecting samples and periodically training the model. For complex models, training phase is typically computationally infeasible which is conflicting with requirements of thesis.

Tesauro et al. [29] 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*. The authors model the data center as multiple applications managed under a single resource manager. Modeling streaming workloads as a queuing system has two problems. First, modeling is a complicated process and determining probability distributions requires domain knowledge. Second, it requires access to each node (so it can be modeled 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. [24] proposed to use Reinforcement Learning to manage resources consumed by virtual machines. It employs standard model-free learning, which is known as *Temporal Difference* [27] 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.

Enda, Enda, and Jim [8] 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. The 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. The 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 from 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).

Heinze et al. [11] implemented Reinforcement Learning in the context of FUGU [9] and compared it to threshold-based approaches. Each node, maintains its own Q-Table and imposes local policy without coordinating other nodes. This architecture can not be applied in the context of spark streaming, since Spark abstracts away individual nodes from the perspective of application developer. In order to decrease state space, the author applied two techniques. First, only system utilization is considered. Second, system utilization is discretized using coarse grained steps. To remedy slow convergence, the controller enforces a *monotonicity constraint* [15]. That is, if the 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.

Cardellini et al. [3] proposed a two level hierarchical architecture for resource management in Apache Storm [26]. There is a local controller on each node which is cooperating with the global controller. The local controller monitors each operator using different policies (threshold-based or Reinforcement Learning using temporal difference). In case, local controller decides to scale in or out an specific operator, it contacts the global controller and informs it about its decision. Then it waits to receive confirmation from the global controller. The global controller operates using a token-bucket-based policy [4] and has global view of cluster. It ranks requests coming from local controllers and either confirms or rejects their decisions. Although, this architecture seems to be a promising approach, however it has been implemented by modifying Storm's internal components. As mentioned above, this is in conflict with thesis's requirements.

In order to mitigate the problem of large state space in Reinforcement Learning, Lolos et al. [22] proposed to start the agent from small number of coarse grained states. As more metrics are collected (and stored as historical records), agent will discover *outlier* parameters (those parameters that are affecting agent more, CPU rather than IO as an example). Then, it partitions the affected state into two states and *re-trains* newly added states using historical records. Both Temporal Difference and Value Iteration methods can be used as learning algorithm. Gradually, agent only focuses on some specific parts of the state space, since all parameters are not equally important. This approach, effectively reduces

the size of state space. However, the trade-off is the storage cost in which historical metrics need to be stored. It worth noting that from the context of paper, storage cost (whether it is in-memory or on-disk and the duration of storing historical metrics) is unclear. Thus, this approach has been abandoned due to uncertainty.

Dutreilh et al. [7] proposed a model-based Reinforcement Learning approach for resource management of cloud applications. All virtual machines are supervised by a single global controller. Slow convergence is the bottleneck of model-free learning, in contrast to model-based learning. However, environment dynamics are not available at the time of modeling. Authors proposed to estimate these parameters as more metrics are collected and then switch to *Value Iteration* [27] algorithm instead of *Temporal Difference*. In short, statistical metrics are stored and updated for each visit of (old state, action, reward, new state) quadruple. As more samples are collected, statistical metrics become more accurate and can be directly used in *Bellman* equation. Until enough measurements get collected, a separate initial reward function is used which is essentially the original reward function but with penalty costs removed. Furthermore, In order to reduce the state space – tuple of [request/sec, number of VMs, average response time] – there exists a predefined upper and lower bound for state variables and average response time is measured at granularity of seconds. This approach has been partially adopted by this thesis.

Dutreilh et al. [6] proposed a model-free Reinforcement Learning approach (*Temporal difference* algorithm) with modified *exploration* policy. The standard exploration policy for Q-Learning is $1 - \epsilon$. Under this policy, the agent performs a random action with probability of ϵ and with probability of $1 - \epsilon$, it adheres to an action proposed by optimal policy. Although the random action is necessary to explore unknown states, but it has severe consequences under streaming workloads. In some cases, it leads to unsafe states where SLOs are severely violated. Since streaming is heavily latency sensitive, this property is undesirable. Thus, author sought toward a heuristic-based policy proposed by Bodik et al. [2]. This policy is based on couple of key observations which has been adopted by this thesis:

- It must quickly explore different states.
- It should collect accurate data as fast as possible, to speedup training.
- During exploration phase, the policy should be careful not to violate SLOs.

8.6 Summary

In this chapter prior work on auto-scaling scaling has been discussed and evaluated. First, threshold-based approaches are investigated. Simple threshold-based approaches are intuitive and simple to understand by application developers and are widely supported by cloud providers. However, adjusting thresholds is a tricky and error-prone process. Then, time-series analysis techniques are explored. As confirmed by other authors, advanced seasonal forecasting is a computationally intensive process, which makes it less suitable for streaming workloads. Queuing theory approaches are suitable for stationary networks with a known probability distribution for workload and service time. Reinforcement Learning techniques has the benefit that it requires zero knowledge about the environment which helps to gradually adapt to changes in environment.

9 Conclusion

Bibliography

- [1] A. Bifet and R. Gavaldà. “Learning from Time-Changing Data with Adaptive Windowing”. In: *Proceedings of the 7th SIAM International Conference on Data Mining*. Vol. 7. Apr. 2007.
- [2] P. Bodik, R. Griffith, C. Sutton, A. Fox, M. I. Jordan, and D. A. Patterson. “Automatic Exploration of Datacenter Performance Regimes”. In: *Proceedings of the 1st Workshop on Automated Control for Datacenters and Clouds*. ACDC ’09. Barcelona, Spain: ACM, 2009, pp. 1–6.
- [3] V. Cardellini, F. L. Presti, M. Nardelli, and G. R. Russo. “Decentralized self-adaptation for elastic Data Stream Processing”. In: *Future Generation Computer Systems* 87 (2018), pp. 171–185.
- [4] V. Cardellini, F. Lo Presti, M. Nardelli, and G. Russo Russo. “Towards Hierarchical Autonomous Control for Elastic Data Stream Processing in the Fog”. In: *Euro-Par 2017: Parallel Processing Workshops*. Ed. by D. B. Heras, L. Bougé, G. Mencagli, E. Jeannot, R. Sakellariou, R. M. Badia, J. G. Barbosa, L. Ricci, S. L. Scott, S. Lankes, and J. Weidendorfer. Cham: Springer International Publishing, 2018, pp. 106–117.
- [5] L. Cherkasova and P. Phaal. “Session-based admission control: a mechanism for peak load management of commercial Web sites”. In: *IEEE Transactions on Computers* 51.6 (2002), pp. 669–685.
- [6] 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.
- [7] X. Dutreilh, S. Kirgizov, O. Melekhova, J. Malenfant, N. Rivierre, and I. Truck. “Using Reinforcement Learning for Autonomic Resource Allocation in Clouds: towards a fully automated workflow”. In: *7th International Conference on Autonomic and Autonomous Systems (ICAS’2011)*. Venice, Italy, May 2011, pp. 67–74. URL: <https://hal-univ-paris8.archives-ouvertes.fr/hal-01122123>.
- [8] 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.
- [9] 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.
- [10] 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.
- [11] T. Heinze, V. Pappalardo, Z. Jerzak, and C. Fetzer. “Auto-scaling techniques for elastic data stream processing”. In: *2014 IEEE 30th International Conference on Data Engineering Workshops*. 2014, pp. 296–302.
- [12] T. Heinze, L. Roediger, A. Meister, Y. Ji, Z. Jerzak, and C. Fetzer. “Online Parameter Optimization for Elastic Data Stream Processing”. In: *Proceedings of the Sixth ACM Symposium on Cloud Computing*. SoCC ’15. Kohala Coast, Hawaii: ACM, 2015, pp. 276–287.
- [13] 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.
- [14] 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.
- [15] 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.

-
- [16] 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.
- [17] M. Kielbowicz. "Online parameter optimization for Spark Streaming". SAP, 2017.
- [18] J. F. C. Kingman. "The Single Server Queue in Heavy Traffic". In: *Proceedings of the Cambridge Philosophical Society* 57 (1961), p. 902.
- [19] 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.
- [20] B. Lohrmann, P. Janacik, and O. Kao. "Elastic Stream Processing with Latency Guarantees". In: *2015 IEEE 35th International Conference on Distributed Computing Systems*. 2015, pp. 399–410.
- [21] B. Lohrmann, D. Warneke, and O. Kao. "Nephele streaming: stream processing under QoS constraints at scale". In: *Cluster Computing* 17.1 (2014), pp. 61–78.
- [22] K. Lolos, I. Konstantinou, V. Kantere, and N. Koziris. "Elastic Resource Management with Adaptive State Space Partitioning of Markov Decision Processes". In: (2017).
- [23] D. A. Menasce, L. W. Dowdy, and V. A. F. Almeida. *Performance by Design: Computer Capacity Planning By Example*. Upper Saddle River, NJ, USA: Prentice Hall PTR, 2004.
- [24] 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.
- [25] 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/.
- [26] A. Storm. *Apache Storm*. Accessed June 21, 2018. 2018. URL: <http://storm.apache.org/>.
- [27] R. S. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction*. 2nd ed. The MIT Press, 1998. ISBN: 0262193981.
- [28] R. Taft, N. El-Sayed, M. Serafini, Y. Lu, A. Aboulmaga, M. Stonebraker, R. Mayerhofer, and F. Andrade. "P-Store: An Elastic Database System with Predictive Provisioning". In: *Proceedings of the 2018 International Conference on Management of Data*. SIGMOD '18. Houston, TX, USA: ACM, 2018, pp. 205–219. ISBN: 978-1-4503-4703-7.
- [29] 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.
- [30] D. Warneke and O. Kao. "Exploiting Dynamic Resource Allocation for Efficient Parallel Data Processing in the Cloud". In: *IEEE Transactions on Parallel and Distributed Systems* 22.6 (2011), pp. 985–997.
- [31] T. Ye and S. Kalyanaraman. "A Recursive Random Search Algorithm for Large-scale Network Parameter Configuration". In: *Proceedings of the 2003 ACM SIGMETRICS International Conference on Measurement and Modeling of Computer Systems*. SIGMETRICS '03. San Diego, CA, USA: ACM, 2003, pp. 196–205.
- [32] Q. Zhang, L. Cherkasova, and E. Smirni. "A Regression-Based Analytic Model for Dynamic Resource Provisioning of Multi-Tier Applications". In: *Fourth International Conference on Autonomic Computing (ICAC'07)*. 2007, pp. 27–27.