

# Auto Scaling of Cloud Resources using Time Series and Machine Learning Prediction

*A thesis submitted for the Degree of  
Master of Science (Research)*

*By*

*Sivasankari Bhagavathiperumal*

*In*

School of Computer Science

UNIVERSITY OF TECHNOLOGY SYDNEY

AUSTRALIA

MARCH 2020

# CERTIFICATE OF ORIGINAL AUTHORSHIP

Date: March 2020

Author: Sivasankari Bhagavathiperumal

Title: Auto Scaling of Cloud Resources Using Time Series and Machine Learning

Degree: Master of Science (Research) in Computing Sciences

I, Sivasankari Bhagavathiperumal declare that this thesis, is submitted in fulfilment of the requirements for the award of Master of Science (Research) in computing Sciences, in the School of Computer Science at the University of Technology Sydney.

This thesis is wholly my own work unless otherwise referenced or acknowledged. In addition, I certify that all information sources and literature used are indicated in the thesis.

This research is supported by the Australian Government Research Training Program.

Signature:      Production Note:  
                         Signature removed prior to publication.

Date: 02/10/2020

# Acknowledgments

My first acknowledgement is to God for being my source of strength. I would like to thank my supervisor Dr. Madhu Goyal for the opportunity to work with her on the thesis, for all the guidance she has given me along the way. I thank the University of Technology Sydney for providing me a wonderful education environment and resources that made this work possible. I would also thank the experts who were involved in validating my project. Finally, I thank my parents, husband and kids for supporting me in my years of study. Thank you all.

# Table of Contents

## Table of Contents

Table of Contents.....	iv
List of Tables .....	vi
List of Figures.....	vii
Abstract.....	1
1. Introduction .....	2
1.1 Significance .....	7
1.2 Motivation .....	8
1.3 Research questions .....	8
1.4 Objective and Aim .....	9
1.5 Organisation of thesis .....	11
2. Literature Review .....	12
2.1 Introduction .....	12
2.2 Auto Scaling .....	14
2.3 Types of Scaling.....	15
2.3 The importance and challenges of auto scaling .....	17
2.4 Approaches of auto scaling .....	19
2.4.1 Linear regression approach .....	19
2.4.2 Spot instances .....	19
2.4.3 Edge computing .....	20
2.4.4 Time series with regression approach.....	21
2.4.5 Deep learning in auto scaling .....	22
2.4.5.1 LSTM and CNN approaches .....	23

2.5 Cost benefits in auto scaling .....	25
2.7 Summary .....	25
3. Framework for Auto Scaling.....	26
3.1 Predictor layer.....	27
3.2 Load receiver.....	29
3.3 Load analyser .....	29
3.4 End user layer .....	30
3.5 Summary of the framework .....	31
4. Data Source.....	32
4.1 Data source details – NASA.....	32
4.2 Data source details – RUET OJ dataset.....	34
4.3 Analysis of datasets.....	36
4.3.1 NASA Dataset analysis.....	37
4.3.2 RUET OJ Dataset analysis .....	39
5. Time Series Prediction .....	41
5.1 ARIMA prediction.....	41
5.2 Exponential smoothing .....	45
6. Machine Learning Prediction .....	48
6.1 Linear regression.....	48
6.2 Naïve method .....	52
6.3 Deep learning (forward propagation) .....	56
7. Conclusion.....	63
7.1 Contribution 1 .....	65
7.2 Contribution 2.....	65
7.3 Contribution 3.....	66
7.4 Future work.....	66
References .....	69

# List of Tables

Table 1: Sample of provisioning resources in Azure.....	14
Table 2:Root mean square error .....	62

# List of Figures

Figure 1: Architecture of cloud-based web applications (Aslanpour, Ghobaei-Arani & Nadjaran Toosi 2017).....	5
Figure 2:Types of Autoscaling .....	16
Figure 4:Proposed Framework load prediction and autoscaling .....	27
Figure 5:Representation of proposed framework.....	28
Figure 6: Flow of data in various layers.....	30
Figure 7:Workload of NASA Dataset (Users Vs Time).....	34
Figure 8:Workload of RUETOJ Dataset (Users Vs Time) .....	35
Figure 9:User Vs Time Vs Date of NASA Dataset .....	37
Figure 10:Time Vs Files Transferred in NASA Dataset .....	38
Figure 11:the load distribution between users and files transferred .....	39
Figure 12:Number of users at particular date and time .....	40
Figure 13: Naïve Forecast for RUET OJ Dataset.....	55
Figure 14:ARIMA prediction for RUET OJ dataset .....	44
Figure 15:Workload forecast using ARIMA of NASA dataset.....	43
Figure 16:Workload prediction using exponential smoothing of NASA dataset .....	47
Figure 17:Sample output of Linear regression for RUET OJ dataset.....	50
Figure 18:Linear Regression method for RUET OJ dataset .....	52
Figure 19:Forecast using Naive approach for NASA Dataset.....	54
Figure 20:Basic layout of neural network .....	57
Figure 21:Basic architecture of feed forward network(Aggarwal 2018) .....	59
Figure 22:Actual and predicted values using forward propagation .....	60

# Abstract

Cloud platforms are becoming very popular as a means to host business applications, and most businesses are starting to use the cloud platform in order to reduce costs. A cloud platform is a shared resource that enables businesses to share the services Software as a service (SAAS), Infrastructure as a service (IAAS) or Anything as a service (XAAS), which are required to develop and deploy any business application. These cloud services are provided as virtual machines (VM) that can handle the end user's requirements. The cloud providers must ensure efficient resource handling mechanisms for different time intervals to avoid wastage of resources. Auto-scaling mechanisms would take care of using these resources appropriately along with providing an excellent quality of service. Therefore, a process to predict the workload is required so that the cloud servers can handle the end user's request and provide the required resources as Virtual Machines (VM) disruptively, so that the businesses using the cloud service only pay for the service they use, thus increasing the popularity of cloud computing. The workload consists of the application programs running on the machine and the users connected to and communicating with the computer's applications. The computing resources are released based on the workload identified using the resource provisioning techniques, which are models used for enabling convenient on-demand network access to a shared pool of configurable computing resources (e.g., networks, servers, applications and services).

In order to provide these resources dynamically without any interruption, an auto-scaling system is required to both manage the load and balance the service. The application should be capable of handling the load as it comes in and to ensure the cloud resources are not supplied in abundance when there are fewer requests coming to the server. This thesis introduces a framework to identify the current workload and predict the future load to the server. The results show the actual workload to the server and the level of requests expected in the future. But to fully appreciate the flexibility of identifying the future load in advance, it was important to identify a method that could in turn provide guidance with adjusting the resources supplied reducing the cost. Also, this would ensure the service level agreement (SLA) was met. The framework identified in this thesis carefully monitors the inputs to the server and analyses the load. The researches also applied deep learning techniques to predict the future load to the server, which produced a result with less root mean square error (RMSE) compared the other techniques. This thesis investigates five different methods to identify the future workload: Naïve, ARIMA, Linear regression, Exponential smoothing and Forward propagation. The logs of two datasets were analysed to predict the future resource for different time intervals along with computing the error between the predicted and actual value.