# Predictive Model for Gap Reduction Between Web Analytics and Business Strategy

by

P.H.A. Nissanka (158230N)

A thesis submitted in partial fulfillment of the requirements for the
Degree of MSc in Computer Science specializing in Cloud Computing

in the
Faculty of Engineering
Department of Computer Science and Engineering

28th February 2019

# Declaration

I declare that this is my own work and contains no material that has been published previously in whole or in part for the fulfillment of any degree program. All the referenced materials have been acknowledged in text.

Student:                                    Supervisor:

_____                    _____
P.H.A. Nissanka (158230N)                   Dr. Shantha Fernando

# *Abstract*

Digital marketing and web analytics are two distinct areas that have captured the attention of many industrial firms. There are a lot of tools developed and a lot of studies carried out in each area separately. But still, a firms ability to harness web analytics to optimize digital marketing elements is limited. This work focuses on evaluating previous work in each of these areas and combine them to build a model that would define the relationship between digital marketing and web analytics. Data captured through each area is expected to be analyzed in the form of a time series forecasting problem. Time series forecasting is a very popular area that captured a lot of firms attention in recent years. This is due to the fact that most real-world problems are linked to a temporal component, and thus can be considered as a time series. Furthermore, this work utilizes cloud services for building and running the learning models.

# Acknowledgements

# Contents

# List of Figures

# List of Tables

# Abbreviations

| | |
|---|---|
| **ANN** | **A**rtificial **N**eural **N**etwork |
| **ARMA** | **A**uto **R**egressive Moving **A**verage |
| **ARIMA** | **A**uto **R**egressive **I**ntegrated Moving **A**verage |
| **AWS** | **A**mazon **W**eb **S**ervices |
| **DeepAR** | **Deep A**uto **R**egressive |
| **GCP** | **G**oogle **C**loud **P**latform |
| **GUI** | **G**raphical User Interface |
| **IaaS** | **I**nfrastructure **as a S**ervice |
| **K-NN** | **K-N**earest **N**eighbour |
| **LSTM** | **L**ong **S**hort-**T**erm Memory |
| **MAD** | **M**ean **A**bsolute **D**eviation |
| **MLMVN** | **M**ulti **L**ayer Multi **V**alued **N**eurons |
| **MSE** | **M**ean **S**quared **E**rror |
| **PaaS** | **P**latform **as a S**ervice |
| **QRF** | **Q**uantile **R**andom **F**orest |
| **RNN** | **R**ecurrent **N**eural **N**etwork |
| **S3** | **S**imple **S**torage **S**ervice |
| **SDK** | **S**oftware **D**evelopment **K**it |
| **TS** | **T**ime **S**eries |