

A Cross Platform Framework for Social Media Information Diffusion Analysis

H.M.M.Caldera

198132D

Doctor of Philosophy

Department of Computer Science and Engineering
Faculty of Engineering

University of Moratuwa
Sri Lanka

November 2023

A Cross Platform Framework for Social Media Information Diffusion Analysis

H.M.M.CALDERA

198132D

Thesis submitted in partial fulfillment of the requirements for the degree
Doctor of Philosophy

Department of Computer Science and Engineering
Faculty of Engineering

University of Moratuwa
Sri Lanka

November 2023

DECLARATION

I declare that this is my own work and this thesis does not incorporate without acknowledgment any material previously submitted for a Degree or Diploma in any other University or institute of higher learning and to the best of my knowledge and belief it does not contain any material previously published or written by another person except where the acknowledgment is made in the text. Also, I hereby grant to the University of Moratuwa the non-exclusive right to reproduce and distribute my thesis, in whole or in part in print, electronic or other medium. I retain the right to use this content in whole or part in future works (such as articles or books).

Signature of the candidate:

Date: 21.11.2023

Mr.H. M. M. Caldera

The above candidate has carried out research for the PhD thesis under my supervision.

Name of the supervisor: Prof. G. I. U. S. Perera

Signature of the supervisor:

.....

Date: 21.11.2023

Acknowledgements

First and foremost, I want to thank my supervisor, Prof. Indika Perera, for his unbeatable support throughout the study. He has taught me, both consciously and unconsciously, how well experimental research is done. I appreciate all his time, ideas, and recommendation for the AHEAD grant to help fund my Ph.D. study.

Next, I would like to thank AHEAD Grants for financial support through research and give special thanks to Dr. Lochandaka Ranathunga as research coordinator for all financial approvals. Dr. Shalinda Adhikari for winning the research grant. I won't complete this degree if you have not won the grant.

As a fellow researcher, Mrs. Nadeera Meedin provided unbelievable support for success during the last few years. She has fully supported my success in many aspects. She shares personal research experiences and motivates me a lot for success.

Further, I thank all the research fellows who worked in the "social media analytics research group" for helping me. Thanks to reviews (including internal ones) for providing valuable feedback.

All authors develop valuable content and publish it on the web. I have gained much knowledge from your videos, documents, research articles, etc. Thanks to all the online content authors. Thanks to Alexander T., Research gate, and other online portals for making many scientific articles publicly available. Further, all scientists who kept their research data as a public asset and let other fellow researchers explore their research.

Finally, I would like to thank my father, mother, wife, son, and all other family members for motivating me and staying with me during this challenging period.

Abstract

In the current digital era, social media platforms have emerged as one of the most effective channels for the diffusion of information. People may readily access and exchange information, news, and opinions from anywhere worldwide because of increasing social media usage.

Information diffusion across multiplex social media platforms is one of the most prominent research problems ever. Social media content generators diffuse information on multiplex social media platforms by targeting many objectives such as popularity, online presence, hate targets, and customer engagement. Regardless of the "content" posted on social media platforms, evaluating the dissemination velocity of each piece of content published on those platforms is essential. It will help to get an overall picture of "how it flows" throughout the social media platforms. Most social media platforms have a platform-specific algorithm for calculating the degree of information diffusion on those platforms. The main objective of this research was to develop a method to calculate the velocity of information diffusion across multiplex social media platforms.

Existing literature on information diffusion strategies, effects, and measurements was used to develop the proposed algorithm. The information diffusion velocity of social media influencers varies according to the content. The platform-specific algorithms for diffusion strength detection vary based on the platform. Somehow, these platform-specific algorithms influence the community to engage with the trending content. i.e., platforms support increasing the strength of information diffusion.

Conventional information diffusion algorithms were designed to measure content diffusion speed on a simplex social media platform, which might be content-specific. The missing dimension is ubiquitous nature. Hence, regardless of the platform, it is mandatory to calculate a ubiquitous information diffusion velocity over multiplex social media platforms.

Both structured information diffusion in a graph for diffusion in a closed network and unstructured patterns in an open-ended coarse-grained information diffusion model check the importance of information diffusion on multiplex social media platforms. Time is another critical factor in defining velocity. i.e., a time series of information diffusion provides a rich picture of information diffusion.

Event-driven architecture is a well-known software architectural approach that facilitates the implementation of microservice-based solutions. The suggested algorithm utilizes an event-driven architecture to manage the information flow by processing social media events. Eventually, this research uses the event-triggering process to understand how information is propagated through an event-driven microservice architecture.

Data science and artificial intelligence are being employed in information diffusion

studies. Understanding how information spreads and the variables and features that influence it is another crucial study area of this research. There are several techniques for studying information dissemination using artificial intelligence. Applying artificial intelligence to information diffusion studies might improve our knowledge of "How information travels" and "how to disseminate information" in various circumstances efficiently. The research used natural language processing to evaluate the textual content of the social media post. That is to find a general textual meaning given by the end-user reactions.

Event-driven architecture is one of the best possible for information diffusion analytics. Using event-driven architecture, data may be delivered in real-time to various analytics services, allowing for the speedy and effective processing of enormous amounts of data. This is especially true in today's data-driven world, when businesses and organizations must make quick, well-informed decisions based on real-time data. Because of its event-driven nature, it is also simple to interface with other systems and services, making it a highly adaptable and versatile option for information distribution analytics. Since the diffusion of information starts with an event's occurrence, it follows numerous steps to flow among the community. An event-driven micro-services architecture that uses artificial intelligence methods (like natural language processing to evaluate textual information) has been experimented with to propose a simple solution for this complex problem.

As per the research work, I can summarize the key findings. I have proposed a tree-structured diffusion tree that can explain how information flows through multiplex social networks. Under this multiplex context, I have experimented with multiple trees and a more robust graph that focused on the diffusion of information. The diffusion strength was based on the SIR model, and the time series analysis focused on how quickly information spread throughout the network. The proposed solution was tested in several real-world cases. Technique-specific tests like seasonality and autocorrelation were conducted to evaluate how the time-series model works in a graph context. Further tests like cohesiveness and robustness were tested, and the proposed algorithm achieved good robustness (an average of 75%) and cohesiveness (an average of 70%) in each case. The best experimental results show an average of more than 80% accuracy in any given instance, and it constructs the tree in less than a second. Most of the predicted values generated an average accuracy of around 70%.

In summary, social media platforms have emerged as prominent channels for information propagation within the contemporary digital landscape. Quantifying the velocity at which information propagates across diverse social networks presents a notable challenge in research. While algorithms tailored to specific platforms influence community engagement, a "universal metric for information dissemination strength" is necessary across multiple social media platforms. The envisioned algorithm considers time series data, integrating structured and unstructured patterns during construction.

Keywords: Information diffusion analysis, Social Media Data Analytics, Graph Learning, Time series analysis, Event-driven micro-services, Artificial Intelligence, Natural Language Processing.

Table of Contents

1	INTRODUCTION	1
1.1	Overview	1
1.2	Problem Statement	2
1.3	Motivation	2
1.4	Research Questions	3
1.5	Objectives	3
1.5.1	Usefulness of the research	3
1.6	Digital social network	4
1.7	Social media networks	5
1.7.1	Social media networks	8
1.7.2	Structure of a Social media network.....	13
1.7.3	Types of social media data.....	13
1.8	Social media analytics (Data science in social media context).....	14
1.8.1	Social trends	14
1.8.2	Data Science	15
1.8.3	Big data.....	15
1.8.4	More on analytical types.....	17
1.9	Getting started with a data science approach.....	18
1.9.1	Data Extraction	18
1.9.2	Feature Engineering.....	19
1.9.3	Feature Engineering for social media data analytics.....	19
1.9.4	Feature Engineering on social media information diffusion	20
1.10	Graph theory and networks	20
1.10.1	Actors and network-level measures.....	21
1.11	Graph data science for social media analytics	22
1.12	Graph data science for measuring information diffusion on network	23
1.12.1	Effect of node level features for information diffusion.....	23
1.13	AI in SM information diffusion context.....	23
1.13.1	Social media information diffusion process	26
1.14	Algorithms used in social media information diffusion process analysis	27
1.15	Time Series analysis	28
1.15.1	Elements of time series data	29
1.15.2	Time series analysis on social media information diffusion	29
1.16	Trend analysis/ Diffusion Analytics.....	29
1.16.1	Techniques used in trend analysis	30

1.16.2	Information diffusion patterns recognition and trend analysis	31
1.17	Simultaneous information diffusion/Cross-Posting	31
1.18	Event-Driven Microservices for the architecture of the system	31
2	Literature Review	34
2.1	Social media networks	34
2.1.1	Characteristics of ties	34
2.1.2	Social Network Sites	35
2.1.3	Propagation of information on social media platforms	36
2.1.4	Software-based Social network analysis	37
2.2	Tree structure and Graph theory in information diffusion context	37
2.2.1	Tree structure for information diffusion analysis	37
2.2.2	Graph theory for information diffusion analysis	38
2.3	Time series analysis	39
2.3.1	Time series techniques	40
2.3.2	Software support for time series analysis	42
2.4	Time series analysis in information diffusion context	43
2.5	Machine Learning algorithms	44
2.5.1	Residual Analysis	46
2.5.2	Trend analysis and forecasting	46
2.5.3	Feature engineering	46
2.5.4	Artificial neural networks	47
2.5.5	Statistical analysis	48
2.5.6	Techniques used in trend analysis	48
2.5.7	Error handling methods	49
2.6	Information diffusion models	49
2.6.1	Herd behavior	50
2.6.2	Graph learning	50
2.6.3	Information Cascade	50
2.6.4	Network-based algorithms used in social media information diffusion process analysis	50
2.6.5	Content-based algorithms for information diffusion analysis	51
2.6.6	Hybrid algorithms	52
2.7	Techniques for Text Preprocessing	52
2.7.1	Summary of Research Areas, Findings, Methods, and References in Social Media Analysis	55
2.8	Event driven architecture	57
2.8.1	Event-driven micro-services	58
2.8.2	Event driven architecture in information diffusion analysis	58
2.9	Model error evaluation	59
3	METHODOLOGY	61
3.1	Data extraction in social media platforms	61
3.1.1	Overview	61
3.2	Research design	62
3.2.1	Solution Overview	63

3.3	Data collection methods	64
3.3.1	YouTube live data extraction	64
3.3.2	Extraction of Data from YouTube Channels in live streaming .	65
3.3.3	Limitations in data extraction	67
3.3.4	Twitter data extraction	68
3.3.5	Facebook data extraction	69
3.3.6	Facebook graphs API	69
3.3.7	Extract using API - general information	70
3.4	Social network analysis	71
3.4.1	ML algorithm overview	71
3.4.2	Contagion approach for information diffusion analysis	71
3.5	Architecture the system for Event-Driven	72
3.5.1	Relationship with the microservices and proposed algorithm .	77
4	ANALYSIS	79
4.1	Proposed Algorithm	79
4.1.1	Overview of the algorithm/Design an algorithm	79
4.2	Adding temporal aspects using time series analysis.	85
4.3	Derivation of an equation	89
4.3.1	Analyze the algorithm behavior	89
4.4	Software and tools for the research	90
4.5	Validity and reliability	91
4.6	Research ethics	93
4.7	Limitations	94
4.7.1	limitations in social media data extraction.	94
5	Evaluation	96
5.1	Data Preparation	96
5.2	Exploratory data analytics	96
5.3	Feature Engineering	100
5.3.1	Feature Extraction	101
5.3.2	Locally Linear Embedding (LLE)	102
5.3.3	Feature Preparation	103
5.3.4	Issues with the number of features	108
5.4	Applying more techniques for fine tuning	109
5.4.1	Feature crossing	109
5.4.2	Hashing	109
5.4.3	Embedding	109
5.5	Defining feature matrix	109
5.6	Statistical analysis	110
5.6.1	Univariate analysis	110
5.6.2	Bi-variate analysis	111
5.7	Statistical forecasting for trend analysis	111
5.7.1	Regression analysis	111
5.7.2	Check for Multicollinearity	119
5.8	Working with Node attributes	119

5.9	Time series analysis.....	122
5.9.1	Statistical techniques.....	123
5.9.2	ARIMA model analysis.....	123
5.9.3	Autoregressive Time Series Modelling.....	123
5.9.4	Techniques of Quantitative Forecasting.....	126
5.9.5	Selecting time series packages.....	127
5.9.6	Detecting missing values.....	128
5.9.7	Evaluating time series data.....	128
5.10	Centrality measures in the Social network.....	129
5.10.1	Edge density distribution.....	131
5.11	Sensitivity analysis.....	133
5.12	contingency table.....	134
5.13	Profiling the micro services.....	134
5.14	Error handling.....	135
5.14.1	Least square method.....	135
5.14.2	Model Evaluation.....	135
5.14.3	Model errors.....	136
5.14.4	Accuracy of algorithm.....	136
6	DISCUSSION	138
6.1	Overview.....	138
7	CONCLUSION	153
7.1	Limitations and drawbacks.....	156
7.2	Future work.....	156
A	Appendix	158
A.1	This code is a sample implementation of the base algorithm.....	158
A.2	This code is a sample implementation of the algorithm that working with a timeseries data.....	159
A.2.1	Parameter definitions for the proposed algorithm.....	162
A.2.2	Diffusion Tree Construction Algorithm Implementation Guide	163
A.3	Introduction.....	163
A.4	Prerequisites.....	163
A.5	Implementation Steps.....	164
A.5.1	Data Preparation.....	164
A.5.2	Algorithm Implementation.....	164
A.5.3	Usage.....	165
A.6	Conclusion.....	166

List of Figures

1.1	Social media network	14
1.2	Drew Conway's Venn diagram of data science	16
3.1	Design Overview	63
3.2	Response JSON Object.....	66
3.3	High-level feature engineering process for data analytical service (Selected based on the key factors	72
3.4	Adopting to contagion approach	73
3.5	High-level overview of the proposed system architecture.....	74
3.6	An overview of proposed RUL situation.....	76
3.7	An overview of proposed event-driven microservices architecture	77
5.1	Box plot for analyzing the outliers.....	99
5.2	Feature selection process	107
5.3	Overview of feature engineering.....	108
5.4	Sample data set	112
5.5	Missing value identification.....	112
5.6	Correlations analysis	113
5.7	Distribution of the number of views	115
5.8	Cumulative distribution of number of views	116
5.9	Distribution of number of likes.....	117
5.10	Implot of number of likes vs. views	119
5.11	category-wise comment distribution.....	120
5.12	ARIMA with Regression	126
5.13	Retweets vs likes in long term diffusion	129
5.14	The network is illustrated using Twitter followers.....	130
5.15	The distribution of edge density	132
5.16	The distribution of edge density, where $n=2$	133
5.17	The distribution of edge density, where $n=2$	134

List of Tables

1	Common features of social media platforms	12
2	Summary of research areas, findings, and methods used in the context of social media analysis.....	55
3	Description of the features.....	97
4	Correlation of the selected attributes	114
5	Regression Model evaluation	118
6	Univariate ARIMA Extrapolation Forecast	124
7	Univariate ARIMA Extrapolation Forecast	125
8	Network base statistical overview.....	131
9	Performance Metrics for Different Algorithms.....	137

List of Abbreviations

<i>AI</i>	Artificial intelligence
<i>API</i>	Application Programming Interfaces
<i>ARIMA</i>	Autoregressive Integrated Moving Average
<i>BC</i>	Betweenness centrality
<i>DLR</i>	Dynamic Linear Regression
<i>EC</i>	Eigenvector centrality
<i>ETS)</i>	Exponential Smoothing
<i>FMTS</i>	Fixed model time series
<i>LDA</i>	Latent Dirichlet Allocation
<i>OMTS</i>	Open model time series
<i>SNA</i>	Social network analysis
<i>SNS</i>	Social Network Sites
<i>STL</i>	Seasonal Decomposition of Time Series
<i>SVM</i>	Support vector machine
<i>TDC</i>	Total degree centrality
<i>TFP</i>	Total-From-Partial