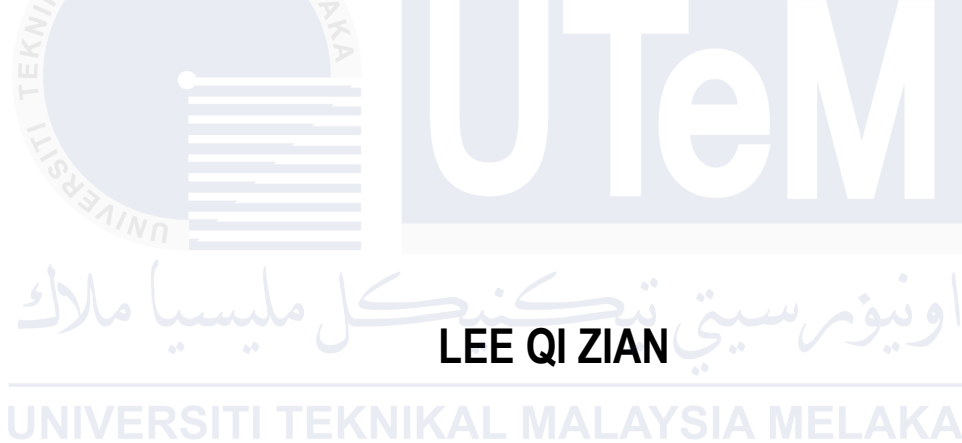




BIG DATA TECHNOLOGY INFORMATION EXTRACTION AND FUSION FROM NON-HOMOGENOUS WEB DATA SOURCES

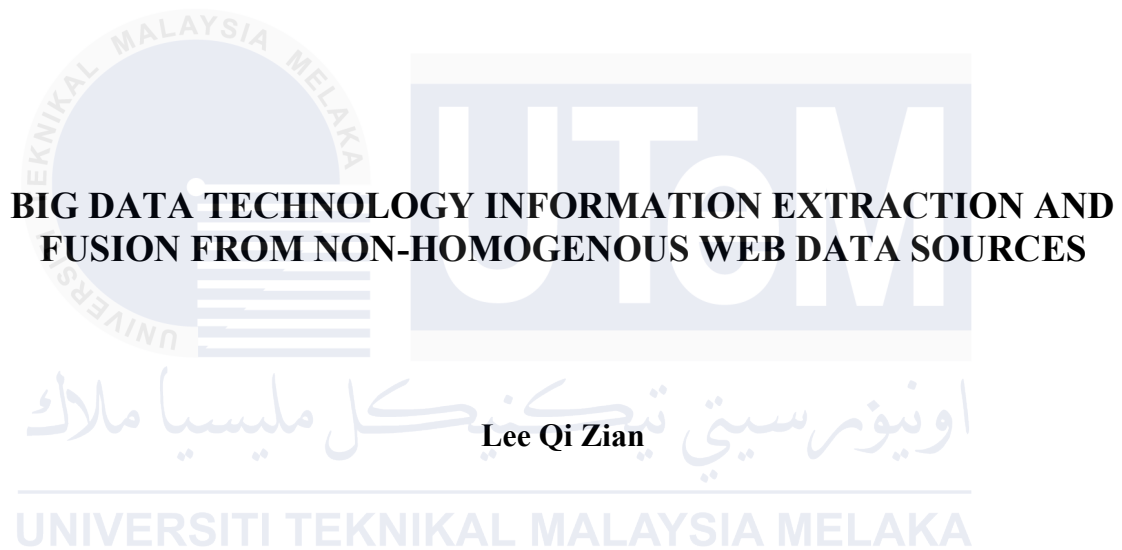


**MASTER OF SCIENCE IN INFORMATION AND
COMMUNICATION TECHNOLOGY**

2025



Faculty of Information and Communication Technology



Master of Science in Information and Communication Technology

2025

**BIG DATA TECHNOLOGY INFORMATION EXTRACTION AND FUSION
FROM NON-HOMOGENOUS WEB DATA SOURCES**

LEE QI ZIAN



**A thesis submitted
in fulfillment of the requirements for the degree of
Master of Science in Information and Communication Technology**



اونيورسيتي تيكنيكل مليسيا ملاك

UNIVERSITI TEKNIKAL MALAYSIA MELAKA

Faculty of Information and Communication Technology

UNIVERSITI TEKNIKAL MALAYSIA MELAKA

2025

DECLARATION

I declare that this thesis entitled “Big Data Technology Information Extraction And Fusion From Non-homogenous Web Data Sources” is the result of my own research except as cited in the references. The thesis has not been accepted for any degree and is not concurrently submitted in candidature of any other degree.



Signature :

Name : LEE QI ZIAN

Date : 11/6/2025

اونيورسيتي تيكنيكل مليسيا ملاك

UNIVERSITI TEKNIKAL MALAYSIA MELAKA

APPROVAL

I hereby declare that I have read this thesis and in my opinion this thesis is sufficient in terms of scope and quality for the award of Master of Science in Information Technology and Communication Technology.



Signature

Supervisor Name

DR NUR ZAREEN ZULKARNAIN

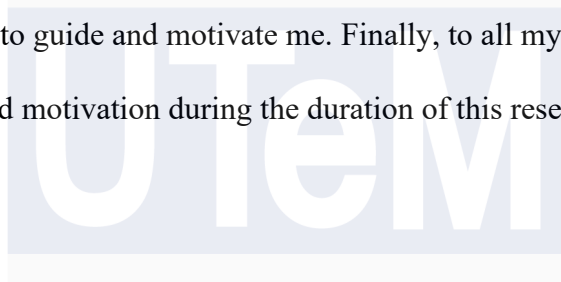
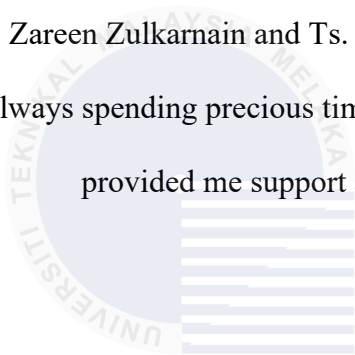
Date

11/6/2025

UNIVERSITI TEKNIKAL MALAYSIA MELAKA

DEDICATION

To my beloved parent Lee Choon Hoe and Yeoh Seok Peng who provide love and support are my greatest inspiration on completing this study. To my dearest supervisor, Dr Nur Zareen Zulkarnain and Ts. Dr Yogan Jaya Kumar for being helpful, responsible and always spending precious time to guide and motivate me. Finally, to all my friends who provided me support and motivation during the duration of this research.



اونيورسيتي تيكنيكل مليسيا ملاك

UNIVERSITI TEKNIKAL MALAYSIA MELAKA

ABSTRACT

Big data has played an ever-increasing role in various sectors of the economy. Despite the availability of big data technologies, many companies and organizations in Malaysia remain reluctant to adopt them. This study was conducted to develop a web extraction framework to extract data from the internet to assist adoption of big data technology. Web scrapping has been a popular method for collecting data from websites. This is because data on the internet is updated frequently thus making it a good source for getting accurate information. Analyzing data requires a large quantity of information to yield a good analysis result. However, the non-homogeneous nature of each website may cause the data from the different internet web sources to have different data making the quality of the data inconsistent. Previous study has propose the use of record linkage method to merge data from multiple website. The record linkage method proposed by previous study used deterministic technique to match data which match the string of matching variable to merge data. However, deterministic technique requires the matching variable to be an exact match to be able to match. Therefore, deterministic matching cannot take into account the dissimilarity such as spacing and different letter cases which can be common in web data due to it non-homogenous nature. This study will explore the use of fuzzy matching technique in matching web data. Fuzzy matching uses Levenshtein distance to calculate the similarity of string and a threshold will be used to decide how similar to trigger a match. This enables fuzzy matching to match string that are only partially match instead of exact match. This study will begin by conducting a systematic review to determine the challenge of big data adoption and what data to extract. This study will implement the Technology-Organization-Environment (TOE) framework to examine the challenges faced by Malaysian organizations with regards to big data adoption. After the systematic review, a web data extraction framework will be developed to extract data that can assist big data adoption. The extracted data will then be merged to enhance the quality of the data. A comparison is made between deterministic matching and fuzzy matching on the performance of merging web data. The finding from this comparison shows that fuzzy matching has a slightly better performance in merging web data. This is due to fuzzy matching can match the string of matching variable that has different spacing and letter cases. A survey case study carried out in this study also shows that the extracted data is very helpful in helping user while purchasing the required big data software on the software market.

PENGEKSTRAKAN MAKLUMAT TEKNOLOGI DATA RAYA DAN GABUNGAN DARIPADA SUMBER DATA WEB YANG TIDAK HOMOGEN

ABSTRAK

Data raya telah memainkan peranan yang semakin meningkat dalam pelbagai sektor ekonomi. Walaupun terdapat teknologi data raya, banyak syarikat dan organisasi di Malaysia masih enggan menerimanya. Kajian ini dijalankan untuk membangunkan rangka kerja pengekstrakan web untuk mengekstrak data daripada internet untuk membantu penggunaan teknologi data raya. Pengikisan web telah menjadi kaedah popular untuk mengumpul data daripada tapak web. Ini kerana data di internet kerap dikemas kini sekali gus menjadikannya sumber yang baik untuk mendapatkan maklumat yang tepat. Proses menganalisis data memerlukan kuantiti maklumat yang banyak untuk menghasilkan keputusan analisis yang baik. Walau bagaimanapun, sifat tidak homogen setiap laman web mungkin menyebabkan data daripada sumber web internet yang berbeza mempunyai data yang berbeza menjadikan kualiti data tidak konsisten. Kajian terdahulu telah mencadangkan penggunaan kaedah pautan rekod untuk menggabungkan data daripada beberapa laman web. Kaedah kaitan rekod yang dicadangkan oleh kajian lepas menggunakan teknik deterministik untuk memadankan data yang sepadan dengan rentetan pembolehubah padanan untuk menggabungkan data. Teknik deterministik memerlukan pembolehubah padanan menjadi padanan tepat untuk dapat dipadankan. Oleh itu, teknik padanan deterministik tidak boleh mengambil kira ketidaksamaan seperti jarak dan kes huruf yang berbeza yang boleh menjadi perkara biasa dalam data web kerana sifatnya tidak homogen. Kajian ini akan meneroka penggunaan teknik logic kabur dalam pemadanan data web. Logik kabur menggunakan jarak Levenshtein untuk mengira persamaan rentetan dan ambang akan digunakan untuk menentukan kesamaan untuk mencetuskan padanan. Ini membolehkan padanan kabur untuk memadankan rentetan yang hanya separa padanan dan bukannya padanan tepat. Kajian ini akan dimulakan dengan menjalankan semakan sistematik untuk menentukan cabaran penggunaan data raya dan data yang perlu diekstrak. Kajian ini akan melaksanakan rangka kerja Technology-Organization-Environment (TOE) untuk mengkaji cabaran yang dihadapi oleh organisasi Malaysia berkaitan penggunaan data raya. Selepas semakan sistematik, rangka kerja pengekstrakan data web akan dibangunkan untuk mengekstrak data yang boleh membantu penerimaan data besar. Data yang diekstrak kemudiannya akan digabungkan untuk meningkatkan kualiti data. Perbandingan dibuat antara padanan deterministik dan logik kabur pada prestasi penggabungan data web. Dapatan daripada perbandingan ini menunjukkan bahawa padanan kabur mempunyai prestasi yang lebih baik sedikit dalam penggabungan data web. Ini disebabkan padanan kabur boleh memadankan rentetan pembolehubah padanan yang mempunyai jarak dan huruf huruf yang berbeza. Walau bagaimanapun, kajian ini juga menunjukkan bahawa data yang diekstrak sangat membantu dalam membantu pengguna semasa membeli perisian data raya yang diperlukan di pasaran perisian.

ACKNOWLEDGEMENTS

First and foremost, I would like to take this opportunity to express my sincere gratitude to my supervisors Dr Nur Zareen Zulkarnain and Ts. Dr Yogan Jaya Kumar from the Faculty of Information and Communication Technology Universiti Teknikal Malaysia Melaka (UTeM) for their essential supervision, support and encouragement towards the completion of this study and thesis.

I would also like to extend my sincere appreciation to the staff of both Universiti Teknikal Malaysia Melaka and Faculty of Information and Communication Technology for assisting me with the administration during my study. I also thank Prof. Dr. Azah Kamilah Binti Draman @ Muda and Centre for Research and Innovation Management (CRIM) for providing financial assistance as part of research assistance allowance during my study.

Finally I would like to extend my gratitude to my family for their unwavering support, love and encouragement, prayers and motivation throughout the duration of this study. Their constant encouragement and support helped me to stay focus and determined to achieve my goals. Without their support, this thesis would not be possible. I am forever grateful to all the person who support me high and lows during this research journey.

TABLE OF CONTENTS

	PAGES
DECLARATION	
APPROVAL	
DEDICATION	
ABSTRACT	i
ABSTRAK	ii
ACKNOWLEDGEMENTS	iii
TABLE OF CONTENTS	iv
LIST OF TABLES	vi
LIST OF FIGURES	vii
LIST OF ABBREVIATIONS	x
LIST OF APPENDICES	xi
LIST OF PUBLICATIONS	xii
 CHAPTER	
1. INTRODUCTION	1
1.1 Background	1
1.2 Problem Statement	2
1.3 Research Objective	5
1.4 Scope of Research	6
1.5 Thesis Outline	7
 2. LITERATURE REVIEW	8
2.1 Introduction	8
2.2 Big Data	8
2.2.1 Data taxonomy	10
2.2.2 Big data tools	12
2.2.1 Big data related works	14
2.3 Data collection	17
2.3.1 Data collection in data science	18
2.4 Web crawling and web scraping	20
2.4.1 Web crawling and web scraping related works	22
2.5 Data fusion in data science	25
2.6 Summary	29
 3. METHODOLOGY	30
3.1 Introduction	30
3.2 Research framework	30
3.2.1 Phase 1: Preliminary study phase	31
3.2.2 Phase 2: Identifying the criteria for big data adoption and the challenges faced by Malaysian companies and organization	32
3.2.3 Phase 3: Exploring data matching techniques for merging data from multiple web sources.	33
3.2.4 Phase 4: Report writing	34

3.3	Summary	35
4.	SYSTEMATIC REVIEW OF BIG DATA CHALLENGES IN MALAYSIA AND DATA COLLECTION	36
4.1	Introduction	36
4.2	TOE Framework	38
4.3	Systematic literature review methodology	40
4.3.1	Selection method	40
4.4	Findings and discussion	41
4.4.1	Challenges from technology context	43
4.4.2	Challenges from organization context	46
4.4.3	Challenges from environment context	48
4.4.4	Comparison of challenges faced by organizations in Malaysia to other countries	50
4.4.5	Mitigation Recommendations	53
4.5	Proposed criteria and data for extraction	55
4.6	Summary	57
5.	WEB DATA EXTRACTION FRAMEWORK AND DATA MERGING	58
5.1	Introduction	58
5.2	Fuzzy matching	60
5.3	Implementation	61
5.3.1	Methodology	62
5.3.2	Determining the similarity threshold for fuzzy matching	66
5.3.3	Evaluating matching validity and comparing data merging performance	67
5.3.4	Result of merging data and comparison between deterministic matching and fuzzy matching	68
5.3.5	Discussion	69
5.4	Case study on information usability of extracted data	71
5.4.1	Survey methodology	71
5.4.2	Discussion of survey findings and analysis	83
5.5	Summary	88
6.	CONCLUSION AND RECOMMENDATIONS	90
6.1	Introduction	90
6.2	Contribution of research	90
6.3	Limitation and future works	92
6.4	Summary	92
	REFERENCES	94
	APPENDICES	113

LIST OF TABLES

TABLE	TITLE	PAGE
2.1	Type of data collected	19
2.2	Data collection tools	19
4.1	Taxonomy of big data adoption challenges pertinent to Malaysian organization	42
4.2	Critical information of software that can assist challenges mitigation.	54
4.3	Data source for extracting data	56
5.1	Additional data merged using fuzzy matching technique compared to deterministic matching technique.	67
5.2	Comparison of number of data merged between deterministic matching technique and fuzzy matching technique.	69
5.3	Survey Demographic	73
5.4	Survey result	76

LIST OF FIGURES

FIGURE	TITLE	PAGE
2.1	Growth of data by year predicted by a study from 2016 to 2025 (Woodie, 2022)	9
2.2	Data Taxonomy	11
2.3	Example of big data implantation in business decision making (Anthopoulos and Kazantzi, 2022)	15
2.4	Data pipeline (Sestino et al., 2020)	16
2.5	Basic design of web crawler (Khalil and Alrub, 2013)	21
2.6	Flow of web crawler extracting data from internet (Baskaran and Ramanujam, 2018)	22
2.7	Web crawler integrated with data pipeline(Baskaran and Ramanujam, 2018)	23
2.8	Web crawler integrated with NLP algorithm from (Kaur, 2022)	24
2.9	Web crawler integrated with sematic reasoning from (Hong et al., 2019)	24
2.10	Statistical matching in combining data from different dataset (Rässler, 2004)	26
2.11	Imputation Scheme (Saporta, 2002)	26
2.12	Matching Variables (Saporta, 2002)	27
2.13	Record linkage method for merging data	27
2.14	Process of deterministic matching	27
2.15	Synthetic data created from merging geo specific census and wealth index (Namazi-Rad et al., 2017)	28
3.1	Research Framework	31
3.2	PRISMA 2020 flowchart	33

3.3	Data Merging (Kim and Park, 2019)	34
4.1	PRISMA flowchart	41
4.2	Challenges of big data adoption in Malaysia	42
4.3	Number of article based on the context of challenges	43
5.1	Example of calculating edits between 2 strings	61
5.2	Process flow of extracting and merging data from web data sources	64
5.3	Data merging using deterministic matching technique.	65
5.4	Data merging using fuzzy matching technique.	65
5.5	Graph of additional data merged using fuzzy matching technique compared to deterministic matching technique.	67
5.6	Dashboard of collected data	73
5.7	Demographic of Respondent Age	74
5.8	Demographic of Respondent Education	75
5.9	Demographic of Industry Sector	75
5.10	Survey result of Description of the software	77
5.11	Survey result of Price	77
5.12	Survey result of Trialability	78
5.13	Survey result of Software Popularity	79
5.14	Survey result of Customer Support	79
5.15	Survey result of Ease of Use	80
5.16	Survey result of Functionality	81
5.17	Survey result of Value for Money	81
5.18	Survey result of Category of user	82
5.19	Survey result of User Insight	83
5.20	Overall finding of survey	84

5.21	Boxplot of respondent data	85
5.22	Correlation matrix of respondent data based on factors	86



LIST OF ABBREVIATIONS

BD	-	Big Data
BDA	-	Big Data Adoption
TOE	-	Technology-Organization-Environment



LIST OF APPENDICES

APPENDIX	TITLE	PAGE
APPENDIX A	ARTICLES BASED ON THE CONTEXT OF CHALLENGES FACE IN BIG DATA ADOPTION IN MALAYSIA	113



LIST OF PUBLICATIONS

The followings are the list of publications related to the work on this thesis:

Qi Zian, Lee and Zulkarnain, Nur Zareen and Jaya Kumar, Yogan, 2024. Challenges in big data adoption for Malaysian organizations: a review. *Indonesian Journal of Electrical Engineering and Computer Science*. 33. 507-517. 10.11591/ijeecs.v33.i1.pp. 507-517. 2024. (SCOPUS indexed, Q2 (2024))

Qi Zian, Lee and Jaya Kumar, Yogan and Zulkarnain, Nur Zareen, 2021. Tools for Big Data and Analytics. *3rd International Conference on Intelligent and Interactive Computing*. pp. 76-79.

CHAPTER 1

INTRODUCTION

1.1 Background

Big data is a technology that had been rapid gaining popularity from organizations looking to seize advantage opportunities. There are three main key characteristic of big data evolution, Volume, Variety, and Velocity (Kaur et al., 2020). Big data can be qualified as any data that are difficult be processed by a traditional data pipeline. Big data has huge application opportunities across all sectors of industries. By adopting big data, organization can leverage their power of its unique characteristic and function to improve the performance of sales and marketing department.

Big data has great potential in research and development due to its unique ability to process data from multiple data sources such as documents, records, protocol, video and images. However, like every technology an organization will always need to invest a significant amount of resources to use it to its fullest potential. The main problem faced by organizations while trying to adopt big data technologies are technology readiness. Despite the simple two words of technology readiness, it has a huge impact on the organization decision on how to move forward. Technologies are an ever-evolving thing and choosing a suitable technology is crucial for the long-term planning of an organization. Decision making has always been in the core of human interaction ranging a simple decision like choosing what to eat during dinner to crucial decision making that can possibly reshape history (Bossaerts and Murawski, 2017). Information has always been a crucial asset in human life since the invention of writing which is also the cornerstone of modern world. By recording

information with words preserves the idea for a long time. They display the vibrations of the human soul in written form, thus the knowledge can be pass down to endless generations. With the invention of computers, human found out that writing can be digitalize and can be store in virtually endless quantities.

Big data has certainly been a tremendous help in the decision making with has largely affect the all the economy sectors (Berardino and Vona, 2023). However, there is no doubt that the financial sector will have the most obvious benefits from big data tools as large volume of data can be process in a much cheaper and efficient ways (Hasan et al., 2020). This can also enable a company to have competitive advantages in the marketing sector such as new product, services, new business model, customer satisfaction and customer loyalties (Ramadan et al., 2020). Moreover, big data can ensure all the financial information of a company to be integrated and explored unlike the traditional model that cause all the information to be kept in silos of each department in the company that can cause crucial information to be missed out (Almeida, 2017).

1.2 Problem Statement

Nowadays there are a lot of tools and software with various capabilities such as Python, Matlab, Tableau, Hadoop, and MongoDB. These tools can provide service in helping organizations and companies tackling big data problems. Although there are a lot of big data tools some companies are still reluctant to use it. In a study conduct in 2018 (Saleh et al., 2018), one of the reason that make Malaysian company especially SME turning away from big data tools is that it can take up a lot of time and resources to use it.

The technologies chosen to tackle big data can usually affect the recruitment of expertise and the training of employee in new technical skills. Choosing the wrong

technologies to use can also affect the quality of the gathered data. According to a study conducted by Gartner in 2017, researchers reveal that poor data quality can cause \$15 million annually and will also undermine their digital initiative which may lead to poor customer trust (Moore, 2018). Despite all the issues, many organizations still underestimate the importance of choosing which tools to be deployed to tackle big data. Business may usually opt for a tool that can be deployed faster, but the tools may be less effective. This issue can also sow distrust between the data execution team and management.

Moreover, the operational cost of big data tools is also a key factor when an organization is deciding on which tools to invest in. A cost of a big data tool ranged from open source (zero cost) to paid license software with good after sales services. Different tools will require different knowledge to operate it such as the programming language and the architecture of the tools. When an organization decides to adopt a new big data tool, they might need to retrain their employees which will result in more cost on the training. Other than operational cost, data sensitivity which is how willingly a company will share their data within department. Data sensitivity has significantly contributed to a factor that makes companies reluctant to adopt big data as one company might have different data sensitivity to other companies (Kaisler et al., 2019).

Even if a company had already decided what type of big data tools to adopt, they will need to find the suitable big data tools on the e-market. The e-market of big data tools tends to be a software review website where software developers or publishers will publish their software on the website for users to review it. The typical review page of a software will usually have the software features, approximate prices, user ratings and the URL link to the software website. These types of websites make finding software easier for end users. However, there are a lot of software tools review websites and all the information on the review

page are varied. For example, a software from review website A may have less listed features compared to another review website B. This can be due to the information in review website A not being updated. One method which can be used to address this issue is by using deterministic matching technique. This technique select a feature in the data as matching variable to match the data thus merging the data sources.

Although previous studies has explored and proven that these implication can be solved by using deterministic matching technique in data merging to match data by comparing the string of the matching variable for match however, deterministic matching can only match data that has exactly same matching variable (Namazi-Rad et al., 2017; Tuoto et al., 2018; Suman et al., 2020; d'Ovidio et al., 2021; Shook-Sa et al., 2022; Jamali-Phiri et al., 2023). The nature of deterministic matching can cause implication in matching data extracted from internet as the data from different sources do not have the exact same matching variables (Shook-Sa et al., 2022).

As a conclusion, there are two major concerns regarding assisting the adoption of big data technology in Malaysia. First is to identify the key challenges concerning Malaysian organizations to use big data technologies. Next is to investigate how relevent information from software tools review websites can be merged to provide useful insights in aiding the adoption of big data technologies. For the latter, the main challenges in merging data from web data source is due to the nature of non-homogenous website that does not have the exact same data. Here we put forward the research questions that will be further investigated in this study.

- What are the key challenges concerning Malaysian organizations to use big data technologies?

- How can we match data regarding big data technologies extracted from different sources for merging without requiring the data to be exactly same for matching?

1.3 Research Objective

To solve and mitigate the problem stated in the problem statement, a study is conducted to develop a framework to assist organizations and companies to adopt big data technologies. The first aim is to find out what are the challenges on big data technologies adoption in Malaysia. This can help us identify what are the helpful information to be extracted from the internet. The second aim is to explore data matching techniques to investigate the optimal techniques in merging data extracted from different data sources. Therefore the main hypothesis of this study is that an alternative data matching techniques such as fuzzy matching are more can match web extracted data better than deterministic matching as it can match data with partially same matching variable. Thus, to achieve the above-mentioned aim, the following objectives have been set:

- a) To identify the challenges concerning big data technology adoption among organization and industries in Malaysia using systematic review.
- b) To propose a web information extraction framework with an optimal approach to extract and combine big data tools information from non-homogenous data sources to assist big data adoption.

To investigate the concerns regarding the use of big data technologies within Malaysian organizations, a systematic review is conducted to identify the criteria and challenges for big data technologies adoption. The information gained from this study is crucial as it will give a clear picture on what type of information are commonly sought by the organization in deciding to use big data technologies.

Next, a web extraction framework will be developed to extract the features concerning of big data tools from the internet based on the information identified from the systematic review. The data extracted from the data sources will be merged. This study will explore an alternative matching technique known as fuzzy matching. The performance of data merging will be compared between deterministic matching and fuzzy matching techniques.

1.4 Scope of Research

The scope of this research are as follows:

- Conducting a systematic review of existing literature to identify the criteria and challenges for the adoption of big data technologies within Malaysian organizations.
- Analyzing the identified criteria and challenges to understand common trends and issues in big data technology adoption.
- Developing a web extraction framework to collect information on features of big data tools from the internet, based on the findings of the systematic review.
- Merging the extracted data from different sources and exploring the use of fuzzy matching as an alternative technique for data merging.
- Comparing the performance of data merging between deterministic matching and fuzzy matching techniques to assess the effectiveness of each method.