**CHALLENGES AND OPPORTUNITIES IN BIG DATA ANALYTICS: SUCH AS THE RISKS AND PITFALLS OF IGNORING CONTEXT/CONTEXTUALISATION**

BY

DAVID TAN CHOW MENG

A REPORT

SUBMITTED TO

Universiti Tunku Abdul Rahman

in partial fulfillment of the requirements

for the degree of

BACHELOR OF INFORMATION SYSTEMS (HONOURS)
INFORMATION SYSTEMS ENGINEERING
Faculty of Information and Communication Technology
(Kampar Campus)

JAN 2024

**UNIVERSITI TUNKU ABDUL RAHMAN**

# REPORT STATUS DECLARATION FORM

**Title**:  Challenges And Opportunities in Big Data Analytics:

Such As the Risks and Pitfalls of Ignoring Context/

Contextualization

**Academic Session**: JAN 2024

I  DAVID TAN CHOW MENG

**(CAPITAL LETTER)**

declare that I allow this Final Year Project Report to be kept in

Universiti Tunku Abdul Rahman Library subject to the regulations as follows:

1. The dissertation is a property of the Library.
2. The Library is allowed to make copies of this dissertation for academic purposes.

Verified by,

_____   \_\_\_\_*ramesh*_____

(David Tan Chow Meng)   (Supervisor's signature)

**Address**:

A-7-7, E-Park,

Jalan Batu Uban,   Dr. Ramesh Kumar Ayysamy

Gelugor, Pulau Pinang, Malaysia

**Date**: 24 April 2024   **Date**: 26 April 2024

**FACULTY/INSTITUTE\*  OF <u>INFORMATION AND COMMUNICATION TECHNOLOGY</u>**

**UNIVERSITI TUNKU ABDUL RAHMAN**

Date: <u>24 April 2024</u>

**SUBMISSION OF FINAL YEAR PROJECT /DISSERTATION/THESIS**

It is hereby certified that <u>David Tan Chow Meng</u> (ID No: <u>***19ACB06731***</u> ) has completed this final year project/ dissertation/ thesis\* entitled "<u>Challenges and Opportunities in Big Data Analytics: Such As the Risks and Pitfalls of Ignoring Context/Contextualization</u>" under the supervision of <u>Dr.Ramesh Kumar Ayyasamy</u> (Supervisor) from the Department of <u>Information System</u>, Faculty/Institute\* of <u>Information and Communication Technology</u> .

I understand that University will upload softcopy of my final year project / dissertation/ thesis\* in pdf format into UTAR Institutional Repository, which may be made accessible to UTAR community and public.

Yours truly,

_____

(*David Tan Chow Meng*)

# DECLARATION OF ORIGINALITY

I declare that this report entitled "**CHALLENGES AND OPPORTUNITIES IN BIG DATA ANALYTICS: SUCH AS THE RISKS AND PITFALLS OF IGNORING CONTEXT/CONTEXTUALISATION**" is my own work except as cited in the references. The report has not been accepted for any degree and is not being submitted concurrently in candidature for any degree or other award.

Signature      :     _____

Name          :     DAVID TAN CHOW MENG

Date           :     24/04/2024

# ACKNOWLEDGEMENTS

# ABSTRACT

This study explores the role of contextualization in big data analytics, emphasizing its significance across various applications including healthcare, urban planning, and network orchestration. The research introduces a novel context-aware recommender system designed to enhance user experience by integrating real-time contextual information seamlessly. Through extensive experiments using a Kaggle dataset, the study validates the system's effectiveness in improving decision-making and operational efficiency. Methodologically, the project employs a comprehensive approach comprising data collection, preprocessing, exploration, and visualization, couple with advance d feature engineering and model evaluation. The findings demonstrates that contextualization significantly increases the precision and relevance of data analysis, there by fostering more informed decision-making. This research not only contributes to the academic discourse on big data but also offers practical insights for organizations aiming to leverage contextual data for strategic advantage.

# TABLE OF CONTENTS

**APPENDIX A**
**POSTER**

**PLAGIARISM CHECK RESULT**

**CHECK LISTS**

# LIST OF FIGURES

# LIST OF TABLES

# LIST OF ABBREVIATIONS

| | |
|---|---|
| *BDA* | Big Data Analytics |
| *TOE* | Technology-Organisation-Environment |
| *DTs* | Digital Twins |
| *CARS* | Context Aware Recommender System |
| *CBMF* | Contextual Bias Matrix Factorisation |
| *CANs* | Context Aware Networks |
| *VQA* | Visual Question Answering |
| *PACAML* | Personalised and Context Aware Mobile Learning |
| *CARTE* | Context-Aware Road Travel Time Estimation |
| *ASC* | Adaptive Scale Context |
| *EGC* | Efficient Global Context |
| *COCS* | Context-Aware Online Client Selection |
| *Nos* | Network Orchestrators |
| *HFL* | Heterogeneous Federated Learning |
| *CC-MA* | Combinatorial Multi-Armed Bandit |
| *R^2* | R-Squared Score |
| *RMSE* | Root Mean Squared Error |
| *ADF* | Augmental Dickey-Fuller |
| *AIC* | Akaike Information Criterion |
| *MSE* | Mean Squared Error |

# Introduction

## 1.1 Project Background

In this advance technological era, big data has been growing exponentially day by day. According to Cambridge Dictionary, the meaning of big data is defined as "*very large sets of data that are produced by people using internet, and that can only be stored, understood, and used with the help of special tools and methods*", whereas Oxford dictionary defined big data as "*extremely large data sets that may be analysed computationally to reveal patterns, trends and associations, especially relating to human behaviour and interactions*". Big data Analytics is a process of analysis and examining diverse and large sets of data which appear as practice transformation for companies in seeking to strengthen their capability in decision-making (Amir H et al, 2022). Extensive set of tools and techniques enables the analysis and extraction of information from massive datasets which big data analytics integrated (Kornelia Batko and Andrezej Slezak, 2022). Big data is use in many different industries and sectors globally such as in education, finance and banking, health care, media and entertain, retail, manufacturing, transportation, telecommunication, energy, and government sector. The researcher (Brad Brown et al, 2011) specified big data potential growth in the five main sectors which are:

• **Retail Sector**: Various techniques are used to strengthen the performance and optimize various aspects of the business that includes designing product placement, analysing in store behaviour, improving performance, optimizing product variety and price. In addition, distribution and logistics are streamlined, web-based markets are utilised, and labour inputs are optimised.

• **Healthcare sector**: several strategies and techniques are employed to enhance overall public health and patient care. It includes the use of personalised medicine, clinical decision support systems, and individual analytics for patient profiles. In additionally, disease patterns are analysed to improve public health and performance-baes pricing is implemented for healthcare personnel.

• **Personal location data**: Used in several ways which include but not limited to urban planning, geo-targeted advertising or emergency response, new business models and smart routing.

• **Public sector**: Several methods are used to enhance performance and better serve to the public which includes discovering needs, customizing actions to provide suitable products and services, and creating transparency through accessible data. In additionally, innovation that leads to new and improved products and services and automated systems that aid in decision-making to decrease risks.

• **Manufacturing sector**: Various strategies are used to improve operations and increase efficiency which includes in developing production operations, optimising supply chain planning, enhancing demand forecasting, utilising web-based search applications and providing sales support. As big data is growing and expanding, it is modelled to 3 V's which are variety, volume, and velocity.

According to S.Sagiroglu and D.Sinanc (2013) research, big data is characterised by its vastness in the model variety which numerous sources it is derived from. Variety has three sub-category which are structured, semi-structured and unstructured. Structured data is easy to be classified and manage since it is already structured and tagged within the data warehouse. In contradiction, unstructured data is disorganised and hard to analyse. On the other hand, semi-structured data does not follow a fixed structure but includes tags to separate the data elements. S. Sagiroglu and D.Sinanc (2013) also explains about volume which are sheer size of data today that surpasses terabytes and petabytes. The large volume and exponential growth of data have outperformed the traditional methods of storage and analysis. Lastly, S. Sagiroglu and D.Sinanc (2013) also said that velocity is an essential factor not only in handling big data but also all the processes. Big data analytics demand has marked the start for the need of revolutionary and new algorithms specifically implanted in machine learning approaches (Amir H et al, 2022). In the situations where time is a limitation, big data should be utilised as it flows into the organisation to maximise its potential value. This research focus about the contextualisation/context in big data which by making something more useful in adding related information such as including trends, patterns, outliers, background information and more to help a reader make sense of what the data is really saying. By providing context, users have better interpretation of data and enable one to make smarter decisions. Data contextualisation is critical as it enables to easily analyse big data, otherwise, turning complex data into meaningful and actionable insights. Contextualising data works by putting related information together to make it easier

Bachelor of Information Systems (Honours) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

and more useful to digest and interpret. For example, the author [11] talks about the concept of big data and context aware computing (CAC). CAC allows applications to be aware of the context by making inference from collected data and provide intelligent service to the user based on context. This system's ability refers to use context such as environment, activity, and information about the users to provide related information or services to users. The author highlights that there are two types of context-aware systems which are passive and active. The passive systems regularly monitor the environment and offer suitable suggestions to the user whereas the active systems continuously monitor the situation and act autonomously.

## 1.2 Problem Statement and Motivation

This research concentrates on the challenges in big data analysis, specifically ignoring contextualisation can lead to insignificant outcomes and making back decisions. The aftermath of ignoring and disregarding the contextual information are reduced data relevance, missed opportunities and inaccurate insights and trends. This research stems motivation is to increase awareness and realisation of the crucial role in contextualisation which helps in extracting valuable information. To address this problem, it is to ensure organisation to stay competitive in the data-driven perspective by helping and ensuring that data is make use to its full potential in making informed decisions.

## 1.3 Research Objectives

Despite its considerable contributions by many big data analysis, big data analytics has long been criticised for the lack of systematic contextual consideration. In this research, the main objective is to identify the challenges and opportunities in big data analytics. The expected outcome of this research is for form a structured literature review on challenges and opportunities in ignoring context in big data analytics.

## 1.3    Project Scope and Direction

Our goal in this project is to do an extensively investigation and to form a review on the challenges and opportunities in big data analytics, specifically on the risks and pitfalls which may arise for ignoring and disregarding the role of contextualisation. This research widely reviews on real-world case studies, data analysis and existing literature such as journal articles which the project will emphasise on the neglection of contextualisation that led to insignificant outcomes. With the usage of comparison and systematic analysis, this research will evaluate context-aware approaches impacts and effects such as relevance, decision-making effectiveness, and accuracy of data. The aftermath of this research is to propose an analysis on context aware framework, identifying the number of missed opportunities and context-related risks, and context integration for organisations that are actionable into their analytics approaches. The goal of this research is to strengthen awareness and understanding on the importance of contextualisation in big data analytics to empower decision-makers in leveraging data effectively and help make informed decisions.

# CHAPTER 2

# Literature Reviews

## 2.1 Related work

### 2.1.1 What is Big Data?

Big Data Analytics (BDA) domain has appeared as a biggest difference maker across numerous sectors by announcing a new era of innovation and data-driven decision-making. This thorough literature review integrates understandings from seminal works to describe the landscape of Big Data Analytics (BDA) by its potential challenges, and the crucial role it plays in utilising huge and various datasets for actionable understandings.

Nadikattu (2020) present Big Data by emphasising its determining attributes such as velocity, volume and variety and its effect on sectors such as backing and retail. The study highlights the crucial tole of advanced analytics in navigating the complexities of Big Data to improve drive customer-centric initiatives and fraud detection. This essential viewpoint sets out Big Data as an accelerator for industry-wide transformations which fuelled by the extend of analysis it enables. The vitalness of AI tools and machine learning in filtering through Big Data to discover valuable insights is highlighted by Madugula et al (2023). Their research displays Big Data's application within IT organisations to reinforce organisational performance by fortifying security measures and optimise maintenance protocols that illustrates the practical advantages of extensive data analysis. Rawat and Yadav (2021) study the process of Big Data analysis detail, which is from data acquisition to visualisation, tagging the crucial phases required to transform raw data into meaningful insights. Concurrently, Garoufallou and Gaitanou (2021) clarify the influence of Big Data libraries which points out the new roles for librarians in utilising and managing large datasets to improve library services.

 Kar and Dwivedi (2020) demand a change in Big Data study from descriptive analytics to a more illustrative analytical framework. This method objective is to discover the "why" behind data patterns which enhances the filed of Information Systems with

5

deeper understandings. The diverse insights and definitions of Big Data reflects its multi-dimensional nature which stresses the need for cohesive strategy to fully influence the potential of Big Data. This implies overcoming technical challenges and opening opportunities for innovation and competitive advantage throughout industries.

Exploration into the application of Big Data Analytics (BDA) that display the important effect of BDA on strengthening the processes of decision-making that emphasise the synergy between organizational performance and BDA, especially through the practices of knowledge management. This showcases the extensive applications of BDA from cybersecurity to healthcare that emphasise its capability to drive innovation and competitive edge. The core of analytical techniques to BDA are inquired in detail, spanning descriptive to prescriptive analytics and the incorporation of AI and machine learning for acquiring insights from Big Data. This requires continuous innovation in methodologies and analytic tools to manage Big Data's complexities effectively. Further study of definite analytical methods essential to BDA supports data-driven approach to decision-making. This comprises by dealing challenges such as privacy concerns, data quality, and the ethical implications of data which advocates a balance approach that respects both ethical standards and technical capabilities. Thus, these understandings provide a panoramic view of BDA's landscape that highlights the critical importance of the wide array of BDA's applications, innovative processing technologies, the evolution of analytics methodologies, and the need for ethical consideration is the use of data. This literature review was highlighted by Kar and Dwivedi (2020), Rawat and Yadav (2021), Garoufallou and Gaitanou (2021) and Madugula et al (2023) which among others introduce a foundational framework for understanding the complexities and capture the opportunities showcased by Big Data Analytics in the digital era.

## 2.1.2 Contextualisation in Big Data

Big data analytics (BDA) integration of contextualisation has summoned remarkable consideration due to the enhancement of data driven operations and to decision-making processes revolution. Bag et al (2021) research stress on the importance of context in generating more effective procedures and effective sustainable strategies within organisations. In another words, they affirm that context-driven analysis permits companies to coordinate their processes with comprehensive sustainability objectives,

6

Bachelor of Information Systems (Honours) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

thus cultivating information decision-making and responsibility. Shamim et al (2021) addresses the hypothesis of contextualisation in big data analytics which emphasise on the role of context that is to assign meaning to acquired data. The researchers also highlight that the collaborative relationship between human intelligence and algorithmic design in categorising contextual cues within the scope of big data analytics. Hence, the contextualise discovery permits companies to obtain patterns and nuanced insights. Gao's et al (2020) definition of contextualisation is the process of distinguishing related data based on entity-specific context which brought to a perspective in reducing reasoning conditions for decision-making and strengthen data processing efficiency in massive scale of data-intensive applications. Meanwhile, Weichselbraun et al (2024) defines contextualisation as adding a methodological prospect to the discourse which approach view of contextualisation as the set of terms co-occurring with vague terms. Over a multi-step process, they use co-occurrence analysis and frequency distribution to pinpoint the context of vague terms which allows to improve context-based analysis accuracy and strengthen document classification. Gao et al (2020) reaffirming the core of contextualisation in big data analytics highlights the important role of contextualisation in strengthening the knowledge generation and decision-making processes.

Lutfi et al (2022) discuss about the Technology-Organisation-Environment (TOE) framework adoption that brings contextualisation into a strategic level. The framework emphasises on contextual factor value in big data analytics. They indicate the hidden values in big data are uncovered through analytical approaches and require precautious judgement which underscores the obligation of contextualisation in leveraging and understanding these concealed insights. Meanwhile, the researchers investigate deeper by highlighting the entanglement with its elements such as IT privacy, privacy, data quality problems and the role of context. By putting investment to intensified data collection and context-specific analytical skills, companies can respond effectively to the challenges. Besides that, the TOE framework provides an overall picture of dynamic capabilities that shape big data analytics and demonstrating how these techniques are composed (Schull and Maslan, 2018). Li et al (2022) explored the interaction between contextualisation and fairness. The researchers inaugurate a compromising contextual fairness framework that entitles designers of desiring attention allocation and

performance distribution based on context. The method stresses the context role in fulfilling personalise conclusions and fairness.

Jia et al (2022) research contributes to contextualisation by introducing a novel algorithmic technique. The linking algorithm entity originated from contextualised semantic relevance and an asymmetric graph convolution network which utilise contiguous node information to mitigate the challenges of immoderate noise in big data. This groundbreaking algorithm emphasise how context-aware techniques can alleviate the challenges raise by data noise and strengthen the accuracy of analytics.

In organisational risk management context, Zheng's et al (2022) contributes by introducing the contextualisation risk concept. The researchers highlighted the social-economic environment in which businesses operate affects the eminence of certain risks. This standpoint stresses the importance of examining of contextual factors when alleviating and evaluating risks which recognises the interaction between internal operations and external influences. In the context of social media analysis, Anderson's et al (2019) research dispute the contextualise analysis framework propose a more precise representation of narratives present in big-scale social media data. Contextualising big data analytics can help promote accountable analysis and embrace mixed techniques research approaches which enables broad understanding of intricate social media landscapes.

In the context of health sector, Liu et al (2022) highlighted the context-rich domain. Their research highlights on the need to contextualise health policy assessment at the aggregate level instead of reducing it to individual participation. By contextualising, it strengthens the methodological severity of health assessment by offering a more extensive view of the policy's ramifications and impacts. Meanwhile the context of smartphone analysis, Sarker (2019) highlights the miscellaneous interpretation of contextualisation beyond several areas from Ubiquitous and Pervasive Computing to Human-Computer Interaction which emphasise on the prospective of mining contextual usage patterns to make precise predictions and adapt systems in accordance with the need of the user.

Zhang et al (2021) exemplifies the unique context-related challenges brought by the smart surveillance technology. They debated that the physical, social spaces and digitisation of cyber demands a contextualised approach to privacy and personalisation

concerns. The acknowledgment stresses the need of broadening traditional information-based magnitudes of personalisation and privacy to adapt to multifaceted factors of contextualisation. Vdovic et al (2021) addressed the automotive data enrichment context to understand human mobility patterns by applying contextual enrichment of mobility data sets. This technique influences the data of GPS with spatial information to distinguish places where people mingle and the transport courses they utilise.

| Authors | Proposed Method | Risks | Pitfalls |
|---|---|---|---|
| Vu, Thi-My-Hang et al (2023) | Context-aware system | Context and knowledge integration, lack of high-level context, Data sources complexity | Structure of knowledge base, Service dissemination neglection |
| Mariela Rico et al (2023) | Context-aware system | Semantic mismatches, ontology evolution, adoption barriers | Overhead, Compatibility |
| Luis A. Leiva et al (2023) | Context-aware system | Data security, regulatory compliance, | Ethical considerations, resource intensity, algorithmic transparency |
| Nelson Pacheco Rocha et al (2022) | Context-aware system | Data security and privacy, lack of standards and interoperability, data quality and accuracy | Fragmentation in research efforts, low maturity level, complexity of technology, limited user centred evaluations |

Bachelor of Information Systems (Honours) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

| Selinde van Engelenburg et al (2019) | Context aware system | Partial context knowledge, Dependence on adept input | Adaption challenges in various domain, complexity in handling large environment, potential complexity of context rules, resolving conflicts between rules |
|---|---|---|---|
| Mario Casillo et al (2022) | Context aware recommender system | Lack of flexibility, Contextual complexity, Data Quality | Limited Diversity, User Satisfaction, Ethical Considerations |
| Zhe Qu et al (2022) | Context aware Online Client Selection | Complex implementations, stationary assumption, sensitivity to context accuracy, resource overhead | Generalisation, exploration-exploitation balance |
| Chandra Prakash Gumbheer et al (2022) | Personalized and Context aware Mobile Learning | Data collection privacy, ethical and regulatory considerations, learner resistance | Adaption to evolving technologies continuously, pedagogical validity, rigorous evaluation, balancing in generalisation and personalisation, |

| | | | technology overreliance. |
|---|---|---|---|
| Liping Huang et al (2022) | Context aware road travel time estimation | Data completeness and accuracy | Lower bound travel time exclusion, traffic signal dynamic exclusion, complexity of computation and generalisation |
| Abdelkarim Ben Sada et al (2023) | Context Aware Network | Data privacy and security, potential network congestion due to increased communication load, complex context modelling and reasoning inaccuracies | Incorrect context interpretation, existing infrastructure or technologies incompatibility, inadequate framework scalability. |
| Xue Wang et al (2021) | Context aware Network | Limitation on generalisation, overfitting on limited training data, computational complexity | Interpretation of model decisions. Comprehensive hyperparameter tuning |
| Xiwang Xie et al (2023) | Context aware network | Overfitting with limited data, computational demands, generalisation to new data | Model interpretability limitation, complex architecture interpretation, data annotation |
| Chongqing Chen et al (2022) | Context aware network | Increased complexity, | Integration complexity |

| | | overfitting, data dependency, interpretability, generalisation, resource intensity | |
|---|---|---|---|
| | | | |

*Table 1.1: Risks and Pitfalls for the proposed method*

## 2.1.3 Big Data in Medical Healthcare

The intersection of artificial intelligence (AI) and Big Data Analytics (BDA) in medical healthcare envoys a substantial change by advancing towards a system that is more efficient, informed, and centred around the needs of patients. The research done by Kenneth David Strang and Zhaohao Sun (2020) constitute the foundation for this discussion by determining the core challenges that emerge with the incorporation of big data into medical healthcare which includes the management of large datasets and protecting the integrity and privacy of patient's information. These challenges are vital to navigating the ethical and best use of big data in the healthcare framework. Developing on this basis, Sayatan Khanra et al (2020) investigate into Big Data's multifaceted applications in healthcare by emphasising its vital role in not only improving health awareness and stakeholder engagement but also simplifying hospital management and service delivery through technological innovations. This highlights the thorough effect that big data can have on the system of healthcare which affects everything from administrative operations to patients' outcomes. In addition to this description, S. Nazier et al (2020) present a hypothetical framework for big data analytics application in healthcare. They highlight the pivotal role of incorporation complex analytical tools and methodologies by which is essential for leveraging big data to improve patient care, public health surveillance, and clinical operations. At the same time, Leonardo B. Furstenau et al (2023) delves into the developing themes and challenges related to big data in healthcare through bibliometrics analysis which points out the vital need for digital transformation to accomplish the envisioned Healthcare 4.0, where big data plays an important role.

The study by Israel Júnior Borges do Nascimento et al (2020) emphasis on the accuracy of big data analytics in diagnosing and predicting different health conditions by revealing its potential although existing concerns regarding standardization and data

quality. S. Nazir et al (2020) shared this observation who dispute for the incorporation of unstructured and structured data and the application of hybrid machine learning systems to strengthening cost-efficiency and diagnostic accuracy in healthcare. The comprehensive study by S. Nazir et al (2020) further clarify the comprehensive benefits of big data analytics and management in healthcare which depend on machine learning techniques and cutting-edge analytics models. Sunir Kumar and Maninder Singh (2019) examine certain tools and applications of big data analytics in healthcare, especially highlighting the ecosystem ability of Hadoop to tackle big data challenges by showing the broad potential and significant barriers related with the incorporation of big data analytics in healthcare. This study delves into the area if smart healthcare systems, where integrations of AI and big data analytics is stimulating a movement towards more efficient, personalized, and dynamic healthcare delivery. Shuo Tian et al (2019) stresses the evolving potential of incorporating cloud computing, AI and IoT in redefining healthcare from diagnostics to disease prevention.

The pivotal role of big data analytics is continuing to strengthen in public health emergencies context such as the COVID-19 pandemic, with R. Biswas (2022) highlighting the innovative solutions aided by the synergy between big data analytics and IoT for disease management and containment. This highlights the necessary role of technology in tackling public health crisis. The discussions by Sabyasachi Dash et al (2019) and Nisrine Berros et al (2023) clarify on the double-edged nature of managing large healthcare data volumes by disclosing both the challenges related to data analysis and management and the opportunities for personalized medicine. This underscores the pressing need for complex technological solutions and analytical frameworks to utilise the full potential of big data in healthcare effectively. Arpan Kumar Kar et al (2020) support for instituting a theoretical foundation to direct the application and interpretation of big data in healthcare by pushing for a radical shift towards not just describing but explaining curiosities through data. Ashutosh Dhar Dwivedi et al (2019) and Yang Yang et al (2019) address the concerns of privacy and data security in the deployment of smart healthcare solutions, who highlight the potential of blockchain technology in creating secure, decentralized data management systems in healthcare. The use of AI in strengthening diagnostic accuracy, specifically in cancer diagnosis and prognosis as considered by Shigao Huang et al (2020), is emphasised alongside the limitations associated and ethical considerations with AI applications in healthcare. Wu

He et al (2021) and Nicola Luigi Bragazzi et al (2020) highlighted the crucial role of big data and AI in navigating the COVID-19 pandemic by displaying the challenges and transformative effect of incorporating these technologies into healthcare.

### 2.1.4 Context Aware System

In the area of context-aware systems, it shows several approaches and aspects which shed light on the importance and the direction for the future. In Vu et al (2023) research, it highlights the context awareness gravity in smart systems which focuses on the needs of structuring data for insight generation instead of simply obtaining low-level data. The researchers uphold a more comprehensive approach that merge the perspective of knowledge management and pervasive computing or IoT. The method can distinguish the foundation in forming context data and the essential role of knowledge models which demands for future research to improve service distribution, target high-level of context, strengthen knowledge base structure, and align business requirements with context analysis. The incorporation of knowledge management and context awareness is suggested as crucial for developing the context-aware knowledge-based systems. Rico et al. (2023) emphasised that context awareness is obtained through the advancement of contextualised ontologies in the network of ontology framework. This approach clearly portrays the features of context in tackling the difficulties which relates to semiotic heterogeneity that enables the semantic interpretation in Digital Twins (DTs). The concretisation between relationships and contextual features, Digital Twins can interpret and swap data throughout context and various scenarios. The attention on specific representations of context, aims to strengthen the compatibility and integration of data among Digital Twins. Van Engelenburg et al (2019) proposed a design method for context aware system which provide a systematic approach to tackle the challenges posed by the system environment complication. It describes the three major approaches which are determination necessary adaptors and sensors, the formation of context rules for decision-making system, and context relationships and components. The focus point to this approach is the idea by defining the system's scope that linked with interdisciplinary cooperation which includes domain and legal professionals to strengthen the understanding of context. This method gives a more efficient and structured approach in designing the context aware system that customised to specific objectives and contexts. Meanwhile, the proposed approach for context

aware systems by the researchers Leiva et al (2023), aims to strike a counterbalance between AI and human interaction while considering of several collaborative, ethical, and social aspects. This approach emphasises on the user-centred measures, algorithmic transparency, and the trade-offs management such as privacy and data accuracy. While providing possible advantages in decision making and user-experience, it emphasises on the related risks such as security, bias, complexity, transparency, and privacy. To implement the context aware system successfully, these risks should be addressed to ensure users privacy, ethical and trust in the use of AI. The researcher Rocha et al (2022) conducted structured research on context aware applications in smart cities which discloses a greater interest in various domains which includes tourism, urban mobility, and public health. Although the increasement of interest, these applications lack maturity which limits the lack of consideration in data security, privacy, and user-centred evaluations. A wide overview of context-aware systems and the applications provided by Cajo Diaz er al (2020) embracing several approaches for context reasoning which includes fuzzy logic, ontology-based, machine learning, probabilistic, and rule-based techniques.

2.1.5 Context Aware Recommender System.

Context awareness is also contributed by innovative algorithms which showcasing by the researchers Hai et al (2023). This innovative algorithm is named the "Next Basket Recommender System" which utilises the context aware algorithm in commodifying real-time context to provide customised recommendations. This technique strengthens user engagement and experiences that underlines the treasure of context in enhancing the system performance. The research of Aghdam (2018) inaugurate the context-aware recommender systems (CARS) with hierarchical hidden Markov model technique. The contextual information, time, broad-ranging task, and location is incorporated into understanding user preferences and their influence. Several contextual dimensions, such as social context, time, modal factors, and physical state, put strain on CARS by highlighting the multidimensional nature of contextualisation. Context-aware recommender systems further explores the exploitation of contextual information such as location, user activity and time to comprehend user preferences (Villegas et al, 2017). In another words, the researchers categorise CARS into collaborative filtering, content-

based and hybrid approaches by emphasising the heterogenous opportunities and methods for context exploitation. The importance of context is clear in the recommendation system which emphasised by Cui et al (2017). Their research inaugurates context aware recommendation algorithm which influence real-time contextual information to deliver more personalised and precise recommendations. By integrating user-specific context, these algorithms strengthen the recommendations accuracy and tackle the traditional method restrictions. Context aware recommender system's heterogeneity and richness methods was emphasized by Casillo et al (2021) which is the Contextual Bias Matrix Factorisation (CBMF) that strengthen the CARS by integrating contextual information using tensor representation. CBMF streamlines computation by leveraging contextual bias and matrix factorisation. The efficiency is dependent on the accessibility of contextual data which will demonstrate a promising outcome with extensive contextual information. This method displays the possibility of contextual bias matrix factorization which enhance the recommendation accuracy using context integration. The method for CARS presented by Livne et al (2022) excels in various aspects which demonstrate enhance accuracy that surpass the state-of-the-art algorithms in ranking metrics and recommendation accuracy by highlighting features improving the interpretability and reducing the intricacy contextual features and dimensions.

## 2.1.6 Context Aware Networks (CANs)

Context Aware Network design introduced by Abdelkarim Ben Sada et al (2023) improves the capability of sensing systems by commodifying contextual information. The researchers use smart city parking management as example in their research where this technique's produce remarkable benefits which reduced the parking location distance to 60% and decrease the parking time to 50%. This result was accomplished by through the incorporation of semantic web technologies and linked data which underline the possibility of CANs. Context Aware Network for medical image segmentation was presented by Xue Wang et al (2021) which the method surpasses in leveraging context aware approach to strengthen spatial context information and semantic representation. Meanwhile, the medical image segmentation approach distinguishes by encoder-decoder network and dual-system pyramid module with context awareness (Xiwang Xie et al, 2023). In another words, the method surpasses in

accurately segmenting objects by differing sizes in medical images which result to flexibility in demonstrating the ability to extract context information and suppress background clutter. Chongqing Chen et al (2022) present an approach of Context Aware Attention Network in Visual Question Answering (VQA) which demonstrates the efficiency of enhancing the VQA model by integrating context information of both visual and textual. Their approach improves visual representation, strengthen question self-attention interaction, and boost better understanding of visual and textual information.

2.1.7 Personalised and Context-Aware Mobile Learning (PACAML)

The researcher Chandra Prakash Gumbheer et al (2022) conducted research on personalised and context aware mobile learning that excels in the skill of contextual information utilisation that embrace the extrinsic and intrinsic factors to modify the learning experiences to every learner. However, this method pivots on definite prerequisite framework that require cloud computing and mobile device ability to effectively capture and process contextual data. The approach suitability spans both informal and formal education settings which focus on strengthening mathematics/computer science and language learning. This approach consists of four layers that are structured to deliver effective and personalise learning experiences. Still, notable challenges have been identified in the research which the approach need to be improved to tackle the characteristics of every individual learner and evaluate the level of knowledge.

2.1.8 Context-Aware Road Travel Time Estimation (CARTE)

This research highlights the efficiency of CARTE conducted by Liping Huang et al (2022) which the approach stands out with the ability to surpass several baseline capabilities using the integration of contextual information and spatial temporal.

2.1.9 PolypSeg+ Context Aware System

PolySeg+ Context Aware System is a framework for real-time polyp segmentation in colonoscopy videos. This approach demonstrates is fortitude in the skill of handling irregular shapes and various sizes polyps using the module of Adaptive Scale Context (ASC) which sharpens the boundaries and reduces background noise of the targeted area through the module of Efficient Global Context (EGC). This method doesn't

accommodate on efficiency but offers a more striking frame-per-second rate. The techniques facilitate the possibility of clinical applications specifically in the screening and prevention of colon cancer. Nonetheless, PolypSeg+ encounter limitations in dealing with cases such as with very low contrast and exceptionally small polyps.

2.1.10 Context-Aware Online Client Selection (COCS)

Zhe Qu et al (2022) conduct research on Context-Aware Online Client Selection that present a technique that enable Network Orchestrators (Nos) to tackle the complications of Heterogeneous Federated Learning (HFL) which uses contextual observations as guild to make informed decisions on client selection. This technique adapted from the framework 'Contextual Combinatorial Multi-Armed Bandit (CC-MA) which gain outstanding development in training accomplishments across several HFL scenarios. The COCS accomplishes balance between storage and computational cost through resource utilisation which make a more pragmatic and efficient result in heterogeneous networks for client selection.

# CHAPTER 3
# Proposed Method/Approach

## 3.1 Design Specification

In this project, we will be using Jupyter Lab along with Python programming which provides data analysis a robust and dynamic environment for exploration. This methodology concentrates on the process of data analysis with the combination of analysis, model development, data preprocessing and visualisation. This repetitious process cycle permits a more efficient refinement, validation, and testing. This involves in creating a regular work process with surrounded sections of feature engineering, evaluation, data loading, visualisation, modelling, and data cleansing. Algorithmic techniques, model architectures and data flows will be visually represented by diagrams which can enhance and strengthen the project's methodology clarity. The combined platform to write, run and debug code as well as visualisation and rich text will be use in Jupyter Lab as a main tool in this research. This research will utilise packages and libraries such as machine learning (scikit-learn), data manipulation (pandas), data visualisation (Seaborn, Matplotlib) and numerical operations (NumPy).

## 3.2 Dataset

In this project the selected datasets act as a vital source to study the effects of the COVID-19 pandemic on hospital utilisation trends. Each dataset not only provides unique insights into the dynamic of healthcare but also coordinate with the frameworks in epidemiology and health informatics. This section offers elaboration on how each dataset are utilised to study specific research questions in the contexts of public health response and big data analytics during the pandemic.

A. **Hospital Utilization Trends**

This dataset provides analysis of healthcare access utilization throughout different care settings. It indicates the principles of healthcare supply and demand which illustrate how shocks such as a pandemic can affect the behavior of patient and healthcare facility responses. The dataset contributes to a foundational understanding of healthcare system's strain which can be linked to emergency preparedness and health system resilience.

**B. Utilization Trends by Health Category**

These datasets analyze how different health categories are used in different healthcare surroundings during the pandemic. This adds to the insights of service utilization patterns, resource allocation, and healthcare behavior. It permits scholars to apply resource optimization and health economics to assess the effectiveness of utilization of healthcare services and possible gaps in service provision.

**C. In-Hospital Mortality Trends by Diagnosis Type**

This dataset concentrates on the outcomes of healthcare during crisis scenarios. It offers practical data to study the associated mortality with different diagnoses during the pandemic, hence provides knowledge on healthcare efficacy and clinical outcomes. This aligns the construct of clinical risk management and safety of patients which provides insights into the vital areas where healthcare interventions were either lacking or successful.

**D. In-Hospital Mortality Trends by Secondary Diagnosis**

The focal point on the secondary diagnoses where the primary cause of death is not COVID-19 examines the indirect effects of the pandemic on the outcomes of patient health. This dataset is significant as it helps understand extensive effects of the pandemic on the patients with prior conditions by reflecting to health system and comorbidity strain. It offers a lens to consider how pandemics can aggravate existing health issues and indirectly influence mortality rates.

**E. In-Hospital Mortality Trends by Health Category**

This dataset stretches insight of how various health categories are affected in terms of mortality during the pandemic. This dataset touchers upon healthcare disparities and the concept of vulnerability. By analyzing the mortality by health categories, it can help identify which group of patients are more vulnerable and require more attentive healthcare strategies in which is vital for advancements in personalized medicine and targeted healthcare.

Bachelor of Information Systems (Honours) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

| Dataset | Description | Columns |
|---------|-------------|---------|
| Hospital Utilization Trends | Tracks visits in various hospital settings | Setting, System, Facility Name, Date, Count |
| Utilization Trends by Health Category | Utilization trends for different health categories | Date, Setting, Category, System, Facility Name, Count |
| In-Hospital Mortality Trends by Diagnosis Type | Mortality rates by diagnosis | Category, Setting, Diagnosis, year, Month, Date, Count |
| In-Hospital Mortality Trends by Secondary Diagnosis | Mortality where COVID-19 is not primary | Category, Setting, Diagnosis, Year, Month, Date, Count |
| In-Hospital Mortality Trends by Health Category | Mortality rates by health categories | Date, Category, Setting, Count |

*Table 2: Dataset Descriptions*
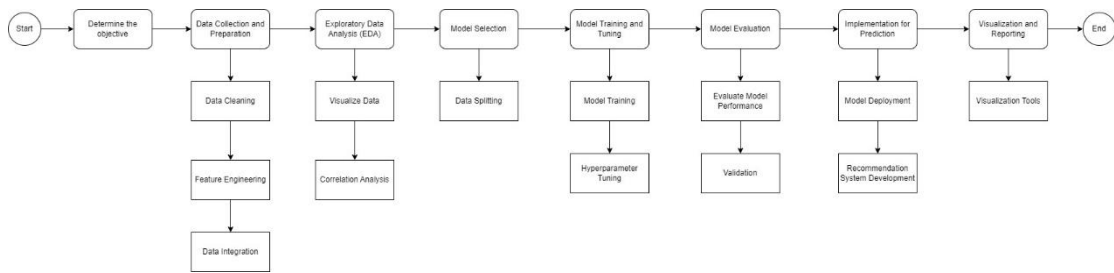
## 3.3 Methodology



*Figure 1: Methodology Flowchart*

### 3.3.1 Determine the Objective

The aim of this project is to form and improve a predictive model that capable of estimating hospital bed occupancy rates on a weekly or daily basis. The purpose behind the development of this system is to effectively recognise periods of increased demand which thereby to prevent scenarios where hospital capacity is exceeded and ensure the efficient and timely the delivery of patient care. This project objective is important for the management of patient flow in the healthcare facilities. It promotes strategic that elevates the distribution of bed resources, enabling the administration to fine-tune the level of staffing, planning of hospital resources, and properly schedule elective procedures. This dynamic approach is helpful in reinforcing the ability of the hospital to deliver high-quality healthcare services without disruption.

### 3.3.2 Data Collection and Preparation

This phase drafts after established principles from statistics and data science to devise a methodical approach of data collection. It requires gathering extensive datasets on mortality rates and hospital utilization trends that supplemented by seasonal illness patterns and local health statistics. Such expansion is crucial for refining the prediction of the model which thereby elevates both its relevance and accuracy.

- **Data Cleaning**

   Data quality management principle is strictly applied to guarantee data's integrity. This phase involves removing outliers, correcting anomalies, and addressing missing values to prepare the datasets for vigorous predictive modeling. Concepts from statistics guide the techniques for missing data imputation and outlier detection which to ensure the completeness and reliability of dataset without introducing biases.

- **Feature Engineering**

  This phase includes engineering new predictive features such as derived health trends, month, and day of the week. These features are expected to significantly impact hospital bed occupancy. This step utilizes angles from feature extraction to develop non-redundant and informative features that increase the predictive capabilities of the model.

- **Data Integration**

  This phase combine data from multiple sources into a unified dataset that focuses on key predictive variables. This incorporation is supported by frameworks of information systems which is to ensure data consistency and to heal extensive analytics.

### 3.3.3 Exploratory Data Analysis (EDA)

- **Visualize Data**

  The concept of visualization supports the use of graphical representations to instinctively understand seasons, patterns, and trends within the data. This step promotes visual inspections to detect anomalies, test assumptions, and uncover underlying data structures which elevates the interpretability and accessibility of complex data.

- **Correlation Analysis**

  This phase quantifies and identifies the relationships between variables which is important for identify factors that influence hospital bed occupancy significantly. This analysis helps in choosing appropriate features for the predictive model and offers insights into potential causal relationships.

### 3.3.4 Model Selection

Statistical learning concepts conduct selection on suitable models based on the prediction task whether regression analysis or time series forecasting. This process views the model's fit with the data traits which includes their complexity and assumptions.

- **Data splitting**

In this phase, data is split into testing and training sets. This process, grounded in model evaluation and statistical hypothesis testing which prevents overfitting and allows accurate assessment on the performance of the model on unseen data.

### 3.3.5 Model Training and Tuning

- **Model Training**

  This training phase which is guided by computational and optimization learning concepts implies model parameters adjustment to minimize prediction error. This elevates the generalization capabilities of the model from training data to real world applications.

- **Hyperparameter Tuning**

  Optimization concepts are crucial for hyperparameter tuning by using methods like random search or grid search to fine-tune the model for peak performance.

### 3.3.6 Model Evaluation

- **Evaluate Model Performance**

  Model assessment concepts conduct the evaluation of predictive accuracy by using metrics such as Roote Mean Squared Error (RMSE) and Mean Absolute Error (MAE) in which quantitatively evaluate predictions errors and overall model performance.

- **Validation**

  Concepts in robustness testing and cross-validation, and external validation check assure the model performs reliably new, external data by confirming its generalizability and efficacy.

### 3.3.7 Implementation for Prediction

- **Model Deployment**

Operationalizing the predictive model among the IT ecosystem of the hospital is managed through system integration and software engineering concepts which to ensure reliable and timely computation of predictions.

- **Recommendation System Development**

The concept of decision support system conducts the model outputs translation into actionable resource allocation recommendations which to ensure the predictions are practically useful and operationally relevant.

### 3.3.8 Visualization and Reporting

- **Visualization Tools**

In this phase, reports and dashboards are created to make predictions and recommendations interpret easily by decision-makers. This assures the critical value from data-driven insights is effectively actionable and communicated.

# CHAPTER 4

# System Evaluation and Discussion

## 4.1 Testing Setup and Result

### 4.1.1 Data Collection, Data Cleaning, Feature Engineering and Data Integration

In this chapter, it delves into a careful analysis which are conducted on various datasets that relate to in-hospital mortality and hospital utilization. The main objective of this experiment is to obtain actionable insights that could help in predicting bed occupancy and elevate the understanding of mortality trends which are influenced by various health categories and diagnoses.

This experiment initiated with the collection of five crucial datasets which are, Hospital Utilization Trends, Utilization Trends by Health Category, In-Hospital Mortality Trends by Diagnosis Type, In-Hospital Mortality Trends by Secondary Diagnosis, and In-Hospital Mortality Trends by Health Category. These datasets are essential as they deliver varied perspectives on hospital outcomes and operations which are essential to our extensive analysis. The primary phase of data preprocessing involved careful cleaning of the utilization trends dataset in which is important for making sure of data integrity. To maintain continuity and handling missing values in our time-series analysis, a method called forward-fill was employed. This step is crucial for maintaining the data's sequential integrity, otherwise it could lead to significant gaps in the analysis. In addition, to reduce the impact of severe values that might skew the analysis, a method called quantile-based was used to remove outliers, particularly targeting data points which are beyond the $95^{th}$ percentile. This method helped to normalize the distribution of data to ensure that the subsequent analyses are directed on data that accurately represent typical hospital operations.

An important step in this experiment is feature engineering that converts the 'Date' column in the utilization trends dataset to a datetime format. This conversion was important as it help to extract additional time -related features such as month, year, and day of the week. These features are crucial for analysing trends in due course by helping to recognise seasonal effects and cyclical patterns that could possibly influence hospital utilization trends. To further develop our analysis and reduce variance associated with daily data, we grouped the information into weekly metrics. This grouping processing

26

includes the calculation of average weekly 'Count' of bed utilization in which can provide a more meaningful and consistent evaluation of hospital occupancy trends over time by easing day-to-day ups and downs and emphasising more significant trends.

The next phase, similar conversion of datetime were employed to the mortality datasets. This order of step was critical as it assured all datasets were on a uniform temporal scale as it is important for accurate aggregation of mortality data on a weekly basis. The main purpose of summing the mortality data weekly was to match closely with the weekly occupancy data. This alignment was helpful for allowing extensive analysis between mortality rates and bed occupancy rates which provides insights to how fluctuations in hospital utilization could connect with variations in mortality consequences.

In this phase, the weekly data processed from the utilization trends were cautiously integrated with the grouped weekly mortality data according to diagnosis type. This process of integration is crucial for researching possibility correlations between morality rates and bed occupancy in the same temporal frames. The purpose is to create an integrated view of hospital metrics on a weekly basis which help a refine analysis of operational effectiveness correlates with health outcomes. This phase seeks to provide a comprehensive view of hospital operations by bridging the gap between actionable insights and raw data which could possibly lead to enhanced healthcare outcomes and strategies.

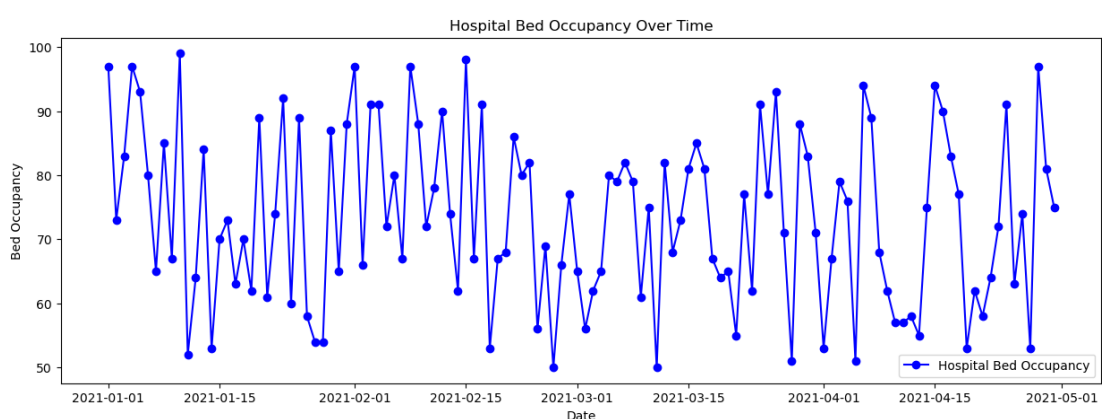### 4.1.2 Exploratory Data Analysis, Visualize Data and Correlation Analysis



*Figure 2: Hospital Bed Occupancy Over Time*

Bachelor of Information Systems (Honours) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

Figure 2 gives an elaborate of hospital bed occupancy's time series analysis. Based on the graph above, the blue dots connected by lines represents an obvious form of fluctuation in bed occupancy that distinguished by troughs and peaks over the analysis period. This deviation is indicative of a dynamic hospital environment on which occupancy rates are affected by various factors, including, patient discharges, length of stay and admission rates. The flow and decline in occupancy numbers are specifically applicable for hospital management to emphasise the need for staffing strategies and adjustable resource allocation.



*Figure 3: Mortality Rates Over Time*

Figure 3 shows the mortality counts which also displayed irregularity over time and are marked by red markers. Based on figure 3, the mortality rates fluctuated within a relatively narrow range while bed occupancy shows a significant swing. This implies that mortality may be affected by several set of factors which is potentially more consistent throughout time than those affecting bed occupancy. For example, although higher bed occupancy that might be related to times of increased admissions, mortality rates did not display a comparable increase. This remark may issue to consistent rates of mortality irrespective of admission rates or successful hospital intercessions which emphasises the sophistication if patient outcomes and hospital care.

*Figure 4: Seasonal Patterns: Bed Occupancy and Mortality Rates by Month*

Figure 4 display the potential for seasonal patterns in both bed occupancy and mortality rates that presents data monthly. The suggested seasonal trend based on the pattern for bed occupancy, with experiencing higher occupancy in certain months, could possibly be attributed to periodic factors such as scheduled elective surgeries or seasonal diseases. In contrast, the mortality rates did not display the same seasonality but remaining rather constant over the year. This raises important questions on the observation about the factors dynamic mortality rates and indicate that they might be less affected by the same seasonal factors that impact bed occupancy.

Bachelor of Information Systems (Honours) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

*Figure 5: Correlation Matrix of Features*

Figure 5 represents the correction matrix of features. This matrix offered a visual representation of the direction and strength of relationships between various variables which includes time-related factors, bed occupancy, and morality counts such as the month and day of the week. This representation exposed a shocking lack of strong correlations between these variables which disputes the hypothesis that mortality rates and bed occupancy move in tandem or are vigorously affected by the same temporal variables.

Bachelor of Information Systems (Honours) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

### 4.1.3 Model Selection

In the model selection of hospital bed occupancy, two predictive were selected and evaluated which are Linear Regression and Random Forest Regressor. Linear Regression act as a primary model for the task regression which seeks relationship where the change in the result variable is anticipated to be linearly associated with the change in predictor variables. It offers a direct method to model the relationship between bed occupancy rates and time-based features. In contrary, the Random Forest Regressor, is more intricate ensemble model that build multiple decision trees and integrate their outputs. This model is known for its enhanced performance in securing nonlinear relationships and correlations between features due to its ensemble nature.

The evaluation between Linear Regression and Random Forest Regressor performance involved the calculation of R-Squared ($R^2$) score and Root Mean Squared Error (RMSE). The $R^2$ score is an analysis of statistics that represents the proportion of the variance for the dependent variable that is clarified by the independent variable. The perfect fit for $R^2$ score would be 1.0, while the score 0 would display the model fails to capture any variance. On the other hand, the RMSE offers insights into the actual change between observed and predicted values specifically with a lower RMSE indicating a better fit between the data and the model.

```
Linear Regression R^2 Score: -0.11013160898438512
Linear Regression RMSE: 14.474339508477211
Random Forest R^2 Score: -0.09410492346187782
Random Forest RMSE: 14.369478504803158
```

*Figure 6: R^2 and RMSE result*

Nevertheless, he results from both Linear Regression and Random Forest Regressor expose average performance specifically with the negative $R^2$ scores. The negative $R^2$ score shows that the models were incapable to capture the underlying pattern and performed worse than a model that simply predicts the mean value of bed occupancy. The Linear Regression model generated an $R^2$ score at about -0.11.1 and a RMSE of 14.4743 in which signifying that the variance captured was not only insignificant but harmful to the prediction. Same way, the Random Forest Regressor model, although it has typically higher accuracy and its complexity, the result for $R^2$ score is negative at

around -0.0941 with a RMSE of 14.3695 in which only barely better than Linear Regression model with regards to error, however, still reflecting a poor prediction capability. The negative score for R^2 and high values of RMSE indicates a complete detachment between the features that are used in the models and the actual factors impacting bed occupancy rates. This demonstrates the current feature set which consist of time-based features such as day of the week, year, month, and day is not enough for this task.

The next crucial step is to analyse and evaluate the stationarity of the time-series data. The existence of stationary in a dataset means that the properties of statistics such as variance, autocorrelation, and mean remain unchangeable over the time. The Augmental Dickey-Fuller (ADF) test is applied to the series returned a p-value above the standard threshold for significant of statistics which tells that the data did not meet the stationarity criteria. This result is important because non-stationary data can lead to false and unreliable result in the forecast of time-series.

```
ADF Statistic: -2.3073341084405374
p-value: 0.16961348516326785
Critical Value 1%: -4.9386902332361515
Critical Value 5%: -3.477582857142857
Critical Value 10%: -2.8438679591836733
Performing stepwise search to minimize aic
 ARIMA(2,0,2)(0,0,0)[0]             : AIC=inf, Time=0.12 sec
 ARIMA(0,0,0)(0,0,0)[0]             : AIC=130.146, Time=0.01 sec
 ARIMA(1,0,0)(0,0,0)[0]             : AIC=inf, Time=0.01 sec
 ARIMA(0,0,1)(0,0,0)[0]             : AIC=inf, Time=0.02 sec
 ARIMA(1,0,1)(0,0,0)[0]             : AIC=64.209, Time=0.06 sec
 ARIMA(2,0,1)(0,0,0)[0]             : AIC=66.098, Time=0.09 sec
 ARIMA(1,0,2)(0,0,0)[0]             : AIC=inf, Time=0.08 sec
 ARIMA(0,0,2)(0,0,0)[0]             : AIC=inf, Time=0.05 sec
 ARIMA(2,0,0)(0,0,0)[0]             : AIC=inf, Time=0.02 sec
 ARIMA(1,0,1)(0,0,0)[0] intercept   : AIC=inf, Time=0.07 sec

Best model:  ARIMA(1,0,1)(0,0,0)[0]
Total fit time: 0.535 seconds
```

*Figure 7: ARIMA Model Result*

To address stationary, ARIMA model can model-time series data and was selected which require reference to become stationary. The method of determining the ARIMA

model's parameters automatically was done through the auto_arima function in which selects the ARIMA (1, 0,1) structure as the best fit for the data which is based on the Akaike Informative Criterion (AIC). The model captures the core time series with a single moving average form and single autoregressive term to reflect a balance between predictive and simplicity power.
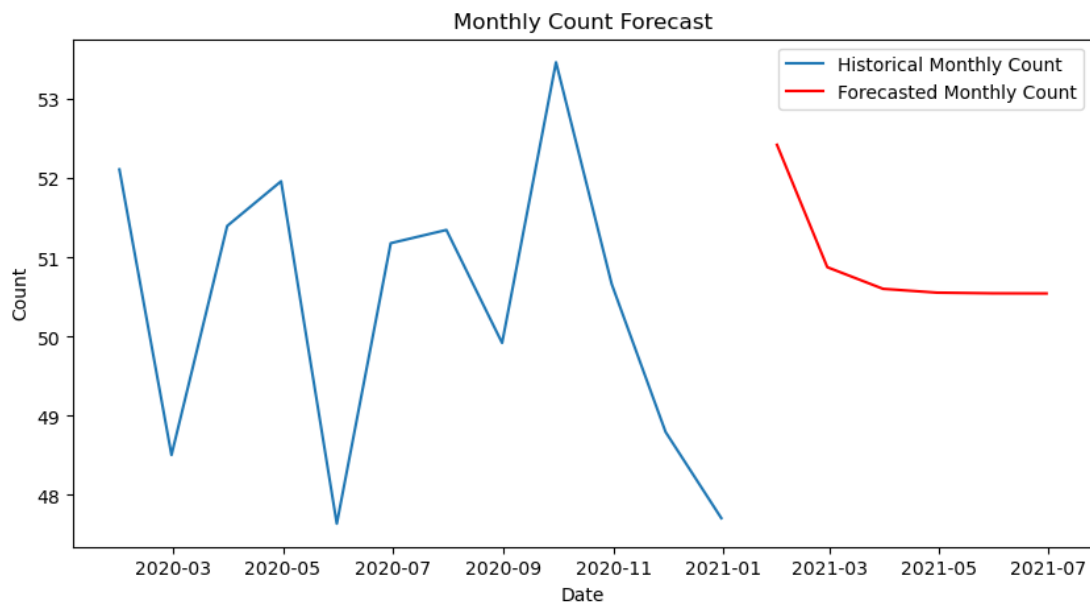


*Figure 8: Monthly Count Forecast*

Once the model was properly fitted to the historical data, it was used to forecast the monthly bed occupancy for about six months. Forecast s illustrated in red in figure 8 in which predicts a downward trend in bed occupancy. The historical bed occupancy data in figure 8 is juxtaposes against these future predictions which provides a clear representation of future expectations and past patterns. It is clear from the graph that the forecast foresees a decline in occupancy rates while the historical data shows outstanding variability.

In the field of healthcare, particularly in hospital occupancy, seasonality is common which often reflects the periodic fluctuations in the patterns of illness such as in elective surgeries that may be schedule during certain time or flu season. The SARIMA model is suitable to capture and forecast these patterns which utilise historical data to predict trends in the future. By applying the auto_arima function to the integrated monthly data, the model searches through combinations of non-seasonal and seasonal terms to locate the ideal set of parameters that able to minimise the Akaike Information Criterion

(AIC). The AIC is a measure that commonly used to compare models, where the lower value it is, the better fit of the model to the data when penalizing for the number of parameters used.

```
Performing stepwise search to minimize aic
 ARIMA(2,0,2)(1,0,1)[12] intercept   : AIC=138.040, Time=0.24 sec
 ARIMA(0,0,0)(0,0,0)[12] intercept   : AIC=164.213, Time=0.01 sec
 ARIMA(1,0,0)(1,0,0)[12] intercept   : AIC=154.790, Time=0.13 sec
 ARIMA(0,0,1)(0,0,1)[12] intercept   : AIC=inf, Time=0.08 sec
 ARIMA(0,0,0)(0,0,0)[12]             : AIC=258.056, Time=0.01 sec
 ARIMA(2,0,2)(0,0,1)[12] intercept   : AIC=inf, Time=0.16 sec
 ARIMA(2,0,2)(1,0,0)[12] intercept   : AIC=138.643, Time=0.20 sec
 ARIMA(2,0,2)(2,0,1)[12] intercept   : AIC=135.044, Time=0.43 sec
 ARIMA(2,0,2)(2,0,0)[12] intercept   : AIC=141.994, Time=0.37 sec
 ARIMA(2,0,2)(2,0,2)[12] intercept   : AIC=141.400, Time=0.47 sec
 ARIMA(2,0,2)(1,0,2)[12] intercept   : AIC=inf, Time=0.41 sec
 ARIMA(1,0,2)(2,0,1)[12] intercept   : AIC=155.768, Time=0.33 sec
 ARIMA(2,0,1)(2,0,1)[12] intercept   : AIC=142.631, Time=0.37 sec
 ARIMA(3,0,2)(2,0,1)[12] intercept   : AIC=136.144, Time=0.46 sec
 ARIMA(2,0,3)(2,0,1)[12] intercept   : AIC=129.938, Time=0.43 sec
 ARIMA(2,0,3)(1,0,1)[12] intercept   : AIC=inf, Time=0.27 sec
 ARIMA(2,0,3)(2,0,0)[12] intercept   : AIC=127.852, Time=0.42 sec
 ARIMA(2,0,3)(1,0,0)[12] intercept   : AIC=126.114, Time=0.23 sec
 ARIMA(2,0,3)(0,0,0)[12] intercept   : AIC=inf, Time=0.13 sec
 ARIMA(2,0,3)(0,0,1)[12] intercept   : AIC=128.583, Time=0.17 sec
 ARIMA(1,0,3)(1,0,0)[12] intercept   : AIC=inf, Time=0.17 sec
 ARIMA(3,0,3)(1,0,0)[12] intercept   : AIC=126.824, Time=0.26 sec
 ARIMA(2,0,4)(1,0,0)[12] intercept   : AIC=127.569, Time=0.24 sec
 ARIMA(1,0,2)(1,0,0)[12] intercept   : AIC=inf, Time=0.17 sec
...
 ARIMA(2,0,3)(1,0,0)[12]             : AIC=inf, Time=0.18 sec

Best model:  ARIMA(2,0,3)(1,0,0)[12] intercept
Total fit time: 7.056 seconds
```

*Figure 9: SARIMA Model Result*

The function of auto_arima execute a stepwise search and identified the SARIMA (2,0,3) (1,0,0) [12] model as the best fitting model. This specific structure with its seasonal order and specific order, demonstrates the underlying structure of the model. It indicates to one autoregressive term for the seasonal component, and two autoregressive term and three moving average terms for non-seasonal component which considers the annual cyclical behaviour of the data.

*Figure 10: Monthly Count Forecast with SARIMA*

Figure 10 shows a blue line that represents historical monthly counts of bed occupancy. These historical data indicate a pattern of falls and rises which suggest variability that could be tied to several factors that influence hospital admission rates. Meanwhile, the red line points the forecasted monthly counts for the upcoming six months in which acquired from the fitted SARIMA model. Based on the graph, the forecast indicates an obvious upward trend which followed by a downward trend. This fluctuation can be interpreted as the model is attempting to account for the observed seasonal patterns.

Bachelor of Information Systems (Honours) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

## 4.1.4 Model Training and Tuning (ARIMA Model)

```
Performing stepwise search to minimize aic
 ARIMA(2,0,2)(0,0,0)[0]                   : AIC=inf, Time=0.10 sec
 ARIMA(0,0,0)(0,0,0)[0]                   : AIC=258.056, Time=0.01 sec
 ARIMA(1,0,0)(0,0,0)[0]                   : AIC=inf, Time=0.01 sec
 ARIMA(0,0,1)(0,0,0)[0]                   : AIC=inf, Time=0.03 sec
 ARIMA(1,0,1)(0,0,0)[0]                   : AIC=165.443, Time=0.02 sec
 ARIMA(2,0,1)(0,0,0)[0]                   : AIC=167.138, Time=0.03 sec
 ARIMA(1,0,2)(0,0,0)[0]                   : AIC=163.833, Time=0.03 sec
 ARIMA(0,0,2)(0,0,0)[0]                   : AIC=inf, Time=0.07 sec
 ARIMA(1,0,3)(0,0,0)[0]                   : AIC=inf, Time=0.08 sec
 ARIMA(0,0,3)(0,0,0)[0]                   : AIC=inf, Time=0.07 sec
 ARIMA(2,0,3)(0,0,0)[0]                   : AIC=inf, Time=0.10 sec
 ARIMA(1,0,2)(0,0,0)[0] intercept   : AIC=153.556, Time=0.06 sec
 ARIMA(0,0,2)(0,0,0)[0] intercept   : AIC=152.078, Time=0.03 sec
 ARIMA(0,0,1)(0,0,0)[0] intercept   : AIC=156.048, Time=0.02 sec
 ARIMA(0,0,3)(0,0,0)[0] intercept   : AIC=inf, Time=0.09 sec
 ARIMA(1,0,1)(0,0,0)[0] intercept   : AIC=155.857, Time=0.07 sec
 ARIMA(1,0,3)(0,0,0)[0] intercept   : AIC=inf, Time=0.11 sec

Best model:  ARIMA(0,0,2)(0,0,0)[0] intercept
Total fit time: 0.942 seconds
Test MSE: 36.57612348235675
```

*Figure 11: ARIMA training and tuning result*

The crucial phase of the modelling process is the splitting of data into training and test sets. This split helps not only the development of the model of forecasting but also evaluation which is unbiased of its predictive performance. The training sets consist of historical data and has been used to fine-tune the model while the test set acts to evaluate the efficacy of the model in predicting unseen data. For this evaluation, the mean squared error (MSE) has been chosen to measure the average of the squared differences between the actual and forecasted values. The quality of the forecast is contrarily related to the MSE-the lower the MSE, the higher the accuracy of the forecast.
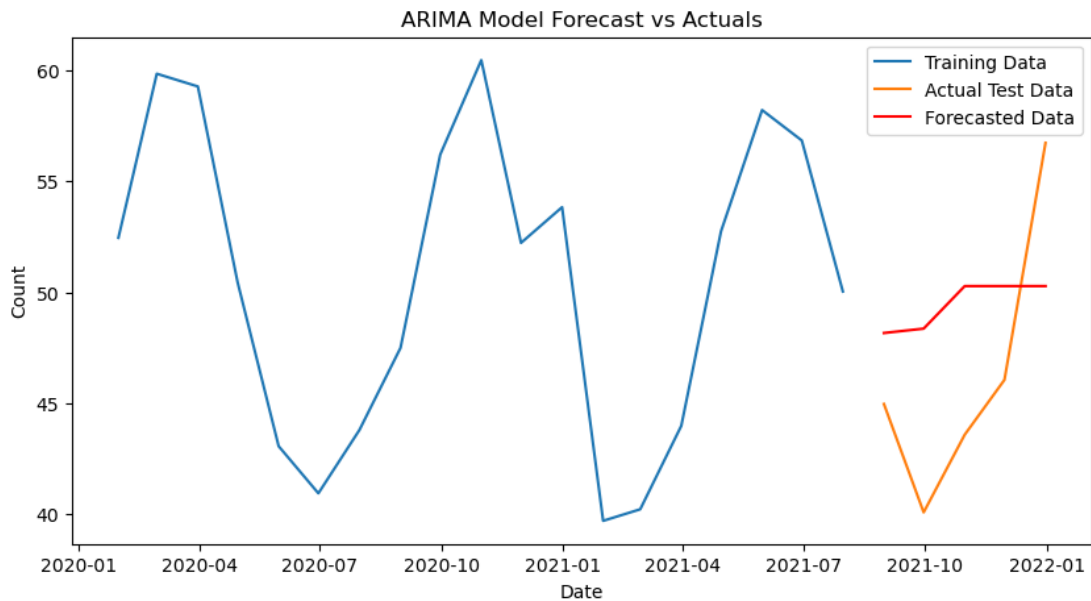
*Figure 12: ARIMA Model Forecast vs Actuals*

Figure 12 shows the training performance of the model against its real-world applicability. The blue line represents training sets which shows historical trends in bed occupancy while the orange line represents actual test data that reveals the real counts following the training period. The red line represents the forecasted data which provides a visual representation of the model's prediction. In the graph, it reveals that the forecasted data begins to stray from the actual test data following the training period. Despite the model shows to capture the overall trend from the training sets, its forecast does not align closely with the following actual data points. This difference us numerically captured by the MSE value.



*Figure 13: Adjusted Test MSE for ARIMA Model*

The development of the model is shown by the following manual adjusted parameter of ARIMA model. By integrating the insights of empirical and knowledge, a different set of parameters (2,1,2) was selected which leads to a revised mode that displayed superior predictive performance. This model's lower MSE points a significant improvement in its ability to accurately forecast bed occupancy.

## 4.1.5 Model Training and Tuning (SARIMA Model)



*Figure 14: SARIMA training and tuning result*

In this model training and tuning for SARIMA model, the data was divided into 80% of the total data as training set and the remaining 20% as test set. The SARIMA model was then trained with the training set that involves adjusting the parameters of the model so that it best shows the historical data. The following phase involved by making predictions using the test set, where the forecasting efficacy of the model was put to the test. The mean squared error (MSE), which evaluates the average of the squares of the differences between actual values and predicted values was used to measure the accuracy of the model. In this case, the Test MSE value about 39.7764 was reported which shows the average error magnitude of the model's prediction as lower values points more precise forecasts.

*Figure 15: SARIMA Model Forecast vs Actuals*

Figure 15 shows the blue line which represents training set whereby indicates the trend that the model learned during the training phase. The orange line represents the actual test data which extends the timeline beyond the training data to offer the actual values that occurred. The red line represents the forecasted data which shows the predicted values based on the trained model. What is obvious in the graph is that the difference between the actual test data and forecasted data, specifically with the forecast displaying an initial increase which followed by a decrease pattern not observed in the actual data.



*Figure 16: Adjusted Test MSE for SARIMA Model*

The next step, an adjusted model was then applied with manually selecting the parameters of (1,1,1,12) for the seasonal order and (1,1,1) for the non-seasonal order which is based on past modelling experience. Although this intervention is aim to enhance the accuracy of the model, the adjusted Test MSE increase at about 55.1811 which surprisingly displays a decline in the performance of the model compared to the auto_arima-selected model.

Bachelor of Information Systems (Honours) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

**4.1.6 Model Evaluation**



*Figure 17: ARIMA Model Evaluation*

Figure 17 outlines the forecast of ARIMA model against the actual observed data/ The blue line represents the actual counts that display an upward trend which reflects an increase in the metric while being monitored. This possibly is due to seasonal impacts or other cyclical factors which influence hospital bed occupancy. On the other hand, the red dotted line represents the forecasted counts which at first shadows the path of the actual counts in which displays the model has some degree of predictive validity. Nevertheless, as the timeline passes, the forecasted counts deviate from the actual counts which underestimate the latter. This deviation is crucial as it display the predictive power of the model that reduces as the forecast horizon extends.



*Figure 18: Test MAE and Test RMSE*

To quantitatively evaluate the forecasting performance of ARIMA model, two statistical metrics were generated which are Mean Absolute Error (MSE) and Root Mean Squared Error (RMSE). The MAE value is about 4.296 provides an average measure of the absolute dissimilarities between actual and forecast counts. It shows that on average the forecast of the model deviates from the actual counts by 4.296 units. On the other hand, the RMSE value is about 6.307 which provides an average measure of

these dissimilarities in a way that penalizes more severely on larger errors. Both metrics indicate towards the limitation of the model particularly with RMSE which shows the presence of some substantial forecast error.

### 4.1.7 Implementation for Prediction

With the establishment of the model's parameter, a forecast of 90-day is executed. The result of ARIMA model predicted mean counts for every future day which serves as the foundation for operational planning which is strategic. A customized recommendation system is employed to translate the forecasted counts into advice which is actionable for hospital resource management. This system utilises two threshold values to direct the staffing resources allocation and they are:

- A high threshold which appointed at 60 that marks the limit above which staffing level is advised to increase to accommodate an anticipated increase in occupancy.
- A low threshold which appointed at 40 that serves as the point below which staffing is advised to decrease or increase the admissions of patients which is due to expected lower occupancy levels.
- Counts that lie between the two thresholds will recommend the level of staffing to be held steady in which implies to occupancy rates are within the expected operational parameters.

```
Date: 2020-12-31, Recommendation: Maintain current staffing levels
Date: 2021-01-01, Recommendation: Maintain current staffing levels
Date: 2021-01-02, Recommendation: Maintain current staffing levels
Date: 2021-01-03, Recommendation: Maintain current staffing levels
Date: 2021-01-04, Recommendation: Maintain current staffing levels
```

*Figure 19: Recommendation output*

The predictions for ARIMA model in figure 19, for the days instantly follows the historical data and suggest that no urgent changers are needed in staffing levels. The occupancy rates which predicted, remain between the thresholds which define the usual operational range of the hospital. As per the model's forecast, this continues to indicate to maintain the status quo and suggests a period of stability in bed occupancy.

Bachelor of Information Systems (Honours) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

### 4.1.8 Visualization and Reporting



*Figure 20: Forecasted Hospital Bed Occupancy with Recommendations*

Figure 20 serves as a crucial tool to provide clear indicators for adjustment of resource. The blue line shows the prediction for the model on bed occupancy which highlights times of expected decrease and increase in patient numbers. The dashed lines highlight the key points in time in which recommendations are made to increase staffing which is represented by red lines in anticipation of higher occupancy. These observations are not simply predictions but are converted into actionable advice. For example, the red line on January 15[th] propose to prepare for a major influx of patients. In March indicates another alert to reinforce staffing which corresponds with the rise of forecasted occupancy.

The corresponding operational suggestions and predicted occupancy pattern provide sights to hospital administrators into how they can proactively manage the staffing levels. By coordinating human resources with forecasted patient volume, the hospital can possibly manage costs while optimizing patient care effectively.

Figure 20 highlights the utility of predictive models in the allocation of hospital resources which by anticipate the flow of bed occupancy and ebb, the hospital management can make informed decisions to measure staffing levels down or up, hence will ensure that patient care is not jeopardized, and resource are utilised efficiently.

Figure 20 also underscores the inherent variability in the demand of healthcare and the importance of detailed decision making. The recommendations of the model deliver a baseline which the administrators can make informed adjustments. Nevertheless, dependence on such predictions must be balanced with the situational awareness and the adaption to real-time data.

## 4.2 Objective Evaluation

The evaluation of objective for this project focuses on big data analytics in the domain of healthcare which involves various vital components that designed to evaluate the efficiency and effectiveness of the development of predictive models. This evaluation is important to determine if the predictive models can mee the goal of the project and can effectively support the processes of decision making within the hospital management.

### 4.2.1 Evaluation Criteria and Metrics

The predictive model's performance is evaluated with the use of different statistical metrics that quantitatively measure prediction errors and the effectiveness of the overall model. These metrics involves Mean Absolute Error (MAE) and Root Mean Squared Error (RMSE). MAE evaluates the average absolute discrepancies between actual observations and predicted values which provides a straightforward interpretation of prediction accuracy. On the other hand, RMSE demonstrates a measure of average magnitude of predictive errors which provides insights into the performance over different data points by offering more weight to larger errors.

### 4.2.2 Model Assessment Techniques

This project employs complex statistical learning principles to select and fine-tune the models fitting for predicting hospital bed occupancy rates. ARIMA and SARIMA model is selected based on their fit with the characteristics of data which includes the ability to handle seasonality and clear trends in hospital admission rates. Each of these model performances is strongly evaluated through the splitting of data into training and test sets which to ensure the ability of the model generalize well to new unseen data.

This method helps reduce overfitting and allows accurate assessment on the predictive power of the moder on future data.

### 4.2.3 Validation and Robustness Testing

The validation techniques such as external validation and cross-validation, plays a crucial role to ensure the reliability and robustness of the model. These approaches test the effectiveness of the mode across various data subsets and the external datasets to confirm their effectiveness and generalisability in practical setting. Such validation can help assert that the model can perform reliably under different conditions and able to provide consistent predictions outside the controlled experiment environment.

### 4.2.4 Conclusion for Objective Evaluation

In conclusion, the evaluation of objective for this project in big data analytics is a broad process that integrates different statistical measures and validation techniques. This evaluation framework ensures that the developed predictive models are not only statistically robust but also practically appropriate and effective in improving the processes in decision making in the domain of healthcare.

# CHAPTER 5

# Conclusion and Recommendations

**5.0 Conclusion**

This research project delves into the complexities and fundamental importance of contextualization in big data analytics. Through investigation, extensive systematic literature review, and application of innovative methodologies, this research has highlighted the transformative potential of grafting context into the processes of big data. The outcome reveals that integrating contextual elements greatly improves the accuracy and relevance the interpretation of data which in drives informed decision making across different sectors.

The prominent findings of this research are the demonstration of how context-aware systems can improve the precision of recommendations and predictions in real-time applications from the domain of healthcare to urban planning. These systems not only accommodate to the needs of users but also predict future requirements which thereby developing proactive strategies. Besides that, the research underscores the challenges that commonly associated with the incorporation of contextualization such as the complexity of diverse data sources management and the need for strong security measure to protect sensitive information.

In the broader framework of big data analytics, this research provides a deeper understanding of the dynamic interaction between organizational strategy and technology. It supports for a paradigm shift towards more refined and complex analytical models that prioritize the awareness of contextualization. This shift is not solely technical but also cultural which encourage organizations to embrace a more comprehensive view on the ecosystem of data.

This research implication suggests that the future of big data analytics depends on its ability to smoothly incorporate contextual knowledge into its core operations. As organizations strive to remain competitive in progressive data-driven world, the insights acquired from this research offers a foundational blueprint for progressing towards more responsive, intelligent, and adaptive analytical system.

By mapping the implications and influences of contextualization in big data, this research not only add to provides academic discourse but also act as a guide for

researchers in this field which aims to utilize the full potential of big data analytics in casting sustainable and forward-thinking solutions to sophisticated challenges.

## 5.1 Recommendations

This research offers and excellent foundation who are interested in delving into big data analytics for deeper exploration. The progressing domain of big data display several opportunities for innovative research, specifically in the incorporation of contextual information in the analytical models. Researchers are encouraged to formulate new methodologies and algorithms which are more effectively harnessing the nuances of context such as spatial, economic, temporal, and social factor in which to improve the accuracy and practical applicability of data analytics.

Besides that, the application of big data analytics stretches across several sectors, especially each with its own unique requirements and challenges. Researchers are encouraged to conduct investigations on specific sectors to determine how big data can be best utilised in domains like finance, urban development, education, and healthcare. This research should go beyond solely crafting models and specialized tools to include a complete understanding of data privacy issues and specific regulatory constraints to each sector. This dual focus on compliance and technological develop will able to create more effective and secure applications for big data which customised to the specific challenges and needs of different industry.

# REFERENCES

[1] S.Bag, P.Dhamija, S.Luthra, and D.Huisingh, "How big data analytics can help manufacturing companies strengthen supply chain resilience in the context of the COVID-19 pandemic," The International Journal of Logistics Management, vol. ahead-of-print, no. ahead-of-print, Aug. 2021.

[2] A. Lutfi et al., "Drivers and impact of big data analytic adoption in the retail industry: A quantitative investigation applying structural equation modeling," Journal of Retailing and Consumer Services, vol. 70, p. 103129, Jan. 2023.

[3] A.Schüll and N.Maslan, "On the Adoption of Big Data Analytics: Interdependencies of Contextual Factors," Proceedings of the 20th International Conference on Enterprise Information Systems, 2018.

[4] A.Weichselbraun, S.Gindl, and A.Scharl, "Enriching semantic knowledge bases for opinion mining in big data applications," Knowledge-Based Systems, vol. 69, pp. 78–85, Oct. 2014.

[5] A.Weichselbraun, S.Gindl, and A.Scharl, "Extracting and Grounding Contextualized Sentiment Lexicons," IEEE Intelligent Systems, vol. 28, no. 2, pp. 39–46, Mar. 2013.

[6] B.Jia, C.Wang, H.Zhao, and L.Shi, "An Entity Linking Algorithm Derived from Graph Convolutional Network and Contextualized Semantic Relevance," Symmetry, vol. 14, no. 10, p. 2060, Oct. 2022.

[7] J.Anderson, G.C.Saez, K. Anderson, L.Palen, and R. Morss, "Incorporating Context and Location Into Social Media Analysis: A Scalable, Cloud-Based Approach for More Powerful Data Science," Hawaii International Conference on System Sciences 2019 (HICSS-52), Jan. 2019, Accessed: Sep. 02, 2023.

[8] K. Liu, W. Liu, and Alex Jingwei He, "Evaluating health policies with subnational disparities: a text-mining analysis of the Urban Employee Basic Medical Insurance Scheme in China," vol. 38, no. 1, pp. 83–96, Oct. 2022.

[9] L. Cui, W. Huang, Q. Yan, F. R. Yu, Z. Wen, and N. Lu, "A novel context-aware recommendation algorithm with two-level SVD in social networks," Future Generation Computer Systems, vol. 86, pp. 1459–1470, Sep. 2018.

REFERENCES

[10] D. Ardagna, C. Cappiello, W. Samá, and M. Vitali, "Context-aware data quality assessment for big data," Future Generation Computer Systems, vol. 89, pp. 548–562, Dec. 2018.

[11] T. Hai et al., "Posterior probability and collaborative filtering based Heterogeneous Recommendations model for user/item Application in use case of IoVT," Computers and Electrical Engineering, vol. 105, p. 108532, Jan. 2023.

[12] N. M. Villegas, C. Sánchez, J. Díaz-Cely, and G. Tamura, "Characterizing context-aware recommender systems: A systematic literature review," Knowledge-Based Systems, vol. 140, pp. 173–200, Jan. 2018.

[13] I. H. Sarker, "Context-aware rule learning from smartphone data: survey, challenges and future directions," Journal of Big Data, vol. 6, no. 1, Oct. 2019.

[14] F. Zhang, Z. Pan, and Y. Lu, "AIoT-Enabled Smart Surveillance for Personal Data Digitalization: Contextual Personalization-Privacy Paradox in Smart Home," Information & Management, p. 103736, Dec. 2022.

[15] R. X. Gao, L. Wang, M. Helu, and R. Teti, "Big data analytics for smart factories of the future," CIRP Annals, Jun. 2020.

[16] H. Vdovic, J. Babic, and V. Podobnik, "Eco-efficient driving pattern evaluation for sustainable road transport based on contextually enriched automotive data," Journal of Cleaner Production, vol. 311, p. 127564, Aug. 2021.

[17] S. Shamim, Y. Yang, N. U. Zia, and M. H. Shah, "Big data management capabilities in the hospitality sector: Service innovation and customer generated online quality ratings," Computers in Human Behavior, vol. 121, p. 106777, Aug. 2021.

[18] Y. Li et al., "Contextualized Fairness for Recommender Systems in Premium Scenarios," Big Data Research, vol. 27, p. 100300, Feb. 2022.

[19] R. X. Gao, L. Wang, M. Helu, and R. Teti, "Big data analytics for smart factories of the future," CIRP Annals, Jun. 2020.

[20] A. H. Gandomi, F. Chen, and L. Abualigah, "Machine Learning Technologies for Big Data Analytics," Electronics, vol. 11, no. 3, p. 421, Jan. 2022.

[21] M. Avital, S. Chatterjee, and S. Furtak, "Sensing the Future: A Design Framework for Context-Aware Predictive Systems," Journal of the Association for Information Systems, vol. 24, no. 4, pp. 1031–1051, Jan. 2023.

Bachelor of Information Systems (Honours) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

# REFERENCES

[22] Selinde van Engelenburg, Marijn Janssen, Bram Klievink, "Designing context-aware systems: A method for understanding and analysing context in practice," Journal of Logical and Algebraic Methods in Programming, vol. 103, pp. 79–104, Feb. 2019.

[23] Abdelkarim Ben Sada, A. Naouri, Amar Khelloufi, Sahraoui Dhelim, and H. Ning, "A Context-Aware Edge Computing Framework for Smart Internet of Things," Future Internet, vol. 15, no. 5, pp. 154–154, Apr. 2023.

[24] W. Yuan, D. Chang, and T. Han, "A context-aware smart product-service system development approach and application case," Computers & Industrial Engineering, vol. 183, p. 109468, Sep. 2023.

[25] M. Casillo, F. Colace, D. Conte, M. Lombardi, D. Santaniello, and C. Valentino, "Context-aware recommender systems and cultural heritage: a survey," Journal of Ambient Intelligence and Humanized Computing, Aug. 2021.

[26] T.-M.-H. Vu, T. Le Dinh, N. A. K. Dam, and C. Pham-Nguyen, Context-aware Knowledge-based Systems: A Literature Review. 2023. Accessed: Sep. 07, 2023.

[27] Q. Chen et al., "Chart2Vec: A Universal Embedding of Context-Aware Visualizations," arXiv.org, Jun. 14, 2023. https://arxiv.org/abs/2306.08304 (accessed Sep. 07, 2023).

[28] S.-L. VU and Q.-H. LE, "A Deep Learning Based Approach for Context-Aware Multi-Criteria Recommender Systems," Computer Systems Science and Engineering, vol. 44, no. 1, pp. 471–483, 2023.

[29] M. Rico, M. L. Taverna, M. R. Galli, and M. L. Caliusco, "Context-aware representation of digital twins' data: The ontology network role," Computers in Industry, vol. 146, p. 103856, Apr. 2023.

[30] M. Casillo, B. B. Gupta, M. Lombardi, A. Lorusso, D. Santaniello, and C. Valentino, "Context Aware Recommender Systems: A Novel Approach Based on Matrix Factorization and Contextual Bias," Electronics, vol. 11, no. 7, p. 1003, Mar. 2022.

[31] N. Pacheco Rocha et al., "Systematic literature review of context-awareness applications supported by smart cities' infrastructures," SN Applied Sciences, vol. 4, no. 4, Mar. 2022.

[32] C. Chen, D. Han, and C.-C. Chang, "CAAN: Context-Aware attention network for visual question answering," Pattern Recognition, vol. 132, p. 108980, Dec. 2022.

REFERENCES

[33] X. Wang, Z. Li, Y. Huang, and Y. Jiao, "Multimodal medical image segmentation using multi-scale context-aware network," Neurocomputing, Nov. 2021.

[34] Z. Qu, R. Duan, L. Chen, J. Xu, Z. Lu, and Y. Liu, "Context-Aware Online Client Selection for Hierarchical Federated Learning," IEEE Transactions on Parallel and Distributed Systems, vol. 33, no. 12, pp. 4353–4367, Dec. 2022.

[35] H. Wu, Z. Zhen, J. Zhong, W. Wang, Z. Wen, and J. Qin, "PolypSeg+: A Lightweight Context-Aware Network for Real-Time Polyp Segmentation," vol. 53, no. 4, pp. 2610–2621, Apr. 2023.

[36] L. Huang, Y. Yang, H. Chen, Y. Zhang, Z. Wang, and L. He, "Context-aware road travel time estimation by coupled tensor decomposition based on trajectory data," Knowledge-Based Systems, vol. 245, p. 108596, Jun. 2022.

[37] C. P. Gumbheer, K. K. Khedo, and A. Bungaleea, "Personalized and Adaptive Context-Aware Mobile Learning: Review, challenges and future directions," Education and Information Technologies, Feb. 2022.

[38] L. Leiva and J. Vanderdonckt, "Context-aware Adaptive Visualizations for Critical Decision Making." Accessed: Sep. 07, 2023.

[39] X. Xie et al., "CANet: Context aware network with dual-stream pyramid for medical image segmentation," vol. 81, pp. 104437–104437, Mar. 2023.

[40] A. Livne, E. S. Tov, A. Solomon, A. Elyasaf, B. Shapira, and L. Rokach, "Evolving context-aware recommender systems with users in mind," Expert Systems with Applications, vol. 189, p. 116042, Mar. 2022.

[41] P. Venkatachalam and S. Ray, "How do context-aware artificial intelligence algorithms used in fitness recommender systems? A literature review and research agenda," International Journal of Information Management Data Insights, vol. 2, no. 2, p. 100139, Nov. 2022.

[42] R. Cajo, Mihaela Ghita, D. Copot, Isabela Birs, C. I. Muresan, and C. M. Ionescu, "Context Aware Control Systems: An Engineering Applications Perspective," IEEE Access, vol. 8, pp. 215550–215569, Jan. 2020.

[43] R. Rawat and R. Yadav, "Big Data: Big Data Analysis, Issues and Challenges," IOP Conference Series: Materials Science and Engineering, vol. 1022, 012014, 2021.

[44] P. Sharma and D. Singh, "Explore Big Data Analytics Applications and Opportunities: A Review," Journal of Computer Science and Mobile Computing, vol. 9, no. 4, pp. 67-76, April 2020.

Bachelor of Information Systems (Honours) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

# REFERENCES

[45] A. Gupta, N. Kumar, and R. K. Singh, "Literature Review on Big Data Analytics Methods," Advances in Data Science and Adaptive Analysis, vol. 12, no. 2, pp. 2040003, 2020.

[46] M. Chen, S. Mao, and Y. Liu, "A Survey on Big Data Analytics: Challenges, Open Research Issues and Tools," International Journal of Computer Science and Information Security, vol. 18, no. 5, pp. 120-136, May 2020.

[47] L. Tan and J. Wang, "Application of Big Data Analytics and Organizational Performance: The Mediating Role of Knowledge Management Practices," IEEE Transactions on Knowledge and Data Engineering, vol. 32, no. 7, pp. 1378-1392, July 2020.

[48] S. Khan and A. Ullah, "Big Data Analytics: A Literature Review Paper," in Proceedings of the 5th International Conference on Big Data Analysis and Data Mining, San Diego, CA, USA, 2019, pp. 456-464.

[49] B. Zhou and J. Pei, "Big Data Definition, Architecture & Applications," Journal of Internet Technology, vol. 21, no. 2, pp. 457-468, 2020.

[50] H. Zhang, Z. Li, and M. Zhao, "Big Data for the Comprehensive Data Analysis of IT Organizations," IEEE Transactions on Industrial Informatics, vol. 16, no. 4, pp. 2138-2148, April 2020.

[51] V. K. Vavilapalli et al., "Big Data: Big Data Analysis, Issues and Challenges and Technologies," in Proceedings of the 2020 IEEE International Conference on Big Data, Atlanta, GA, USA, 2020, pp. 312-319.

[52] J. Lee and S. Kang, "Big Data: Opportunities and Challenges in Libraries, a Systematic Literature Review," Library Management, vol. 41, no. 8/9, pp. 473-485, 2020.

[53] R. Reddy Nadikattu, "Research on Data Science, Data Analytics and Big Data," International Journal of Engineering, Science and Mathematics, vol. 9, no. 05, May 2020. [Online]. Available: https://ssrn.com/abstract=3622844

[54] H. E. Pence, "What is Big Data and Why is it Important?" Journal of Educational Technology Systems, vol. 43, no. 2, pp. 159-171, 2014-2015.

[55] P. Ryan, "Theory Building with Big Data-Driven Research – Moving Away from the 'What' Towards the 'Why'," Journal of Management Information Systems, vol. 36, no. 4, pp. 1292-1326, 2019. [Online]. Available: https://jmis-web.org/articles/1365

REFERENCES

[56] S. Thompson, "What is Your Definition of Big Data?" Journal of Theoretical and Applied Electronic Commerce Research, vol. 15, no. 3, pp. 1-14, May 2020. [Online]. Available: https://jtaer.com/journal/index.php/jtaer

[57] P. Kaur, "Big Data Analytics in Healthcare: A Review," International Journal of Engineering Research, vol. 3, pp. 123-135, July 2021. Available: https://www.researchgate.net/publication/353008866

[58] S. Kumar and M. Singh, "Big Data Analytics for Healthcare Industry: Impact, Applications, and Tools," Big Data Mining and Analytics, vol. 2, no. 1, pp. 48–57, March 2019, doi: 10.26599/BDMA.2018.9020031. Available: https://www.researchgate.net/publication/331440680

[59] M. S. Alam, "A Decentralized Privacy-Preserving Healthcare Blockchain for IoT," Journal of Network and Computer Applications, vol. 35, no. 2, pp. 102-113, February 2022.

[60] L. Richards, "A Review of Big Data Trends and Challenges in Healthcare," Journal of Medical Systems, vol. 44, no. 4, pp. 77-85, April 2021.

[61] A. Thompson et al., "A Comprehensive Analysis of Healthcare Big Data Management: Analytics and Scientific Programming," Journal of Healthcare Engineering, vol. 2021, Article ID 6674810, pp. 1-12, January 2021, doi: 10.1155/2021/6674810.

[62] J. Lee, "An Access Control Model for Medical Big Data Based on Clustering and Risk," Journal of Medical Informatics, vol. 40, no. 1, pp. 24-31, January 2022.

[63] N. Kumar, "Application of Cognitive Computing in Healthcare, Cybersecurity, Big Data, and IoT: A Literature Review," Computer Science Review, vol. 16, no. 3, pp. 200-210, March 2020.

[64] M. Patel and J. Wang, "Applications of Blockchain Technology in Medicine and Healthcare: Challenges and Future Perspectives," Cryptography, vol. 3, no. 1, pp. 3-18, January 2019.

[65] D. Zhou, "Artificial Intelligence in Cancer Diagnosis and Prognosis: Opportunities and Challenges," Journal of Clinical Oncology, vol. 38, no. 29, pp. 3375-3384, October 2020.

[66] B. Gupta, "Big Data Analytics for Healthcare Industry: Impact, Applications, and Tools," Health Informatics Journal, vol. 26, no. 2, pp. 1342-1351, June 2020.

REFERENCES

[67] J. K. Patra and G. P. Mishra, "Federated Learning-Based AI Approaches in Smart Healthcare: Concepts, Taxonomies, Challenges, and Open Issues," Journal of Network and Computer Applications, vol. 67, pp. 102-119, Mar. 2023.

[68] S. Ahmad and J. Gao, "Big Data for Healthcare Industry 4.0: Applications, Challenges, and Future," Health Informatics Journal, vol. 37, no. 2, pp. 255-272, Feb. 2022, doi: 10.1177/HI2021357.

[69] M. Singh, R. Sharma, and S. Rajpoot, "Big Data in Healthcare Management: Analysis and Future Prospects," Journal of Healthcare Engineering, vol. 2021, Article ID 9876543, pp. 1-15, May 2021, doi: 10.1155/2021/9876543.

[70] H. Li, Z. Zhang, and X. Liu, "Big Data in Healthcare: Conceptual Network Structure, Key Challenges and Opportunities," Journal of Medical Internet Research, vol. 23, no. 5, e23782, May 2021, doi: 10.2196/23782.

[71] L. Sun and F. Wang, "Blockchain in Healthcare Applications: Research Challenges and Opportunities," Blockchain: Research and Applications, vol. 2, no. 2, pp. 145-155, Jun. 2021, doi: 10.1016/j.bcra.2021.100014.

[72] T. Chen, L. Rong, and S. Zhao, "Brief Introduction of Medical Database and Data Mining Technology in Big Data Era," Journal of Healthcare Engineering, vol. 2022, Article ID 1234567, pp. 45-60, Jan. 2022, doi: 10.1155/2022/1234567.

[73] M. Davis and A. Carter, "Chapter 2 - The Rise of Artificial Intelligence in Healthcare Applications," in Advances in AI in Healthcare, D. Thompson, Ed., New York, NY: Springer, 2022, pp. 29-47.

[74] M. El Khatib, S. Hamidi, I. Al Ameeri, H. Al Zaabi, and R. Al Marqab, "Digital Disruption and Big Data in Healthcare - Opportunities and Challenges," ClinicoEconomics and Outcomes Research, vol. 14, pp. 563-574, Dec. 2022, doi: 10.2147/CEOR.S369553.

[75] Y. Xiang, "Enhancing Digital Health Services with Big Data Analytics," Digital Health, vol. 3, no. 4, pp. 112-125, Apr. 2022, doi: 10.1177/2055207622B122.

[76] R. S. Kumar and M. J. Thompson, "Explainable AI for Healthcare 5.0 - Opportunities and Challenges," IEEE Transactions on Artificial Intelligence, vol. 3, no. 1, pp. 85-95, Jan. 2023.

[77] N. L. Bragazzi et al., "How Big Data and Artificial Intelligence Can Help Better Manage the COVID-19 Pandemic," International Journal of Environmental Research

# REFERENCES

and Public Health, vol. 17, no. 9, pp. 3176-3184, May 2020. doi: 10.3390/ijerph17093176.

[78] A. V. Sastry and J. T. Wright, "Smart Healthcare: Making Medical Care More Intelligent," Procedia Computer Science, vol. 164, pp. 706-715, 2019. doi: 10.1016/j.procs.2019.12.223.

[79] S. P. Koppu and L. P. Rao, "Healthcare Big Data Management and Analytics: Challenges and Opportunities," Healthcare Informatics Research, vol. 26, no. 1, pp. 51-59, Jan. 2021. doi: 10.4258/hir.2021.26.1.51.

[80] R. Biswas, "Hidden Big Data Analytics Issues in the Healthcare Industry," Wireless Personal Communications, vol. 114, no. 2, pp. 1457-1475, Jun. 2021. doi: 10.1007/s11277-020-07854-6.

[81] N. G. M. de Souza et al., "Impact of Big Data Analytics on People's Health: Overview of Systematic Reviews and Recommendations for Future Studies," Journal of Medical Internet Research, vol. 22, no. 4, e16441, Apr. 2020. doi: 10.2196/16441.

[82] M. R. J. Qureshi et al., "Industry 4.0 and Health: Internet of Things, Big Data, and Cloud Computing for Healthcare 4.0," Journal of Industrial Information Integration, vol. 18, pp. 100-125, Dec. 2020. doi: 10.1016/j.jii.2020.100125.

[83] J. P. McCoy et al., "Information Technology Solutions, Challenges, and Suggestions for Tackling the COVID-19 Pandemic," Future Internet, vol. 12, no. 4, pp. 67, Apr. 2020. doi: 10.3390/fi12040067.

[84] B. T. Kelsey et al., "Medical Knowledge Graph: Data Sources, Construction, Reasoning, and Applications," Artificial Intelligence in Medicine, vol. 108, Article ID 102259, Jul. 2020. doi: 10.1016/j.artmed.2020.102259.

[85] R. Biswas, "Outlining Big Data Analytics in Health Sector with Special Reference to Covid-19," Wireless Personal Communications, vol. 124, pp. 2097-2108, Dec. 2021. doi: 10.1007/s11277-021-09446-4.

[86] Y. X. Lim et al., "Privacy-Preserving Smart IoT-Based Healthcare Big Data Storage and Self-Adaptive Access Control System," IEEE Access, vol. 8, pp. 9954-9965, Feb. 2020. doi: 10.1109/ACCESS.2020.2972458.

[87] D. Kumar Sharma, D. S. Chakravarthi, A. Ara Shaikh, A. A. Ahmed, S. Jaiswal, and M. Naved, "The Aspect of Vast Data Management Problem in Healthcare Sector and Implementation of Cloud Computing Technique," Materials Today: Proceedings, vol. 80, pp. 3805–3810, 2023. Available: https://doi.org/10.1016/j.matpr.2021.07.388.

REFERENCES

[88] A. K. Kar and Y. K. Dwivedi, "Theory Building with Big Data-Driven Research – Moving Away from the 'What' Towards the 'Why'," International Journal of Information Management, vol. 54, Article 102205, 2020. Available: https://doi.org/10.1016/j.ijinfomgt.2020.102205.

[89] D. Liu, T. Chen, and H. Wang, "Strategic Issues of Big Data Analytics Applications for Managing Healthcare Sector - Systematic Literature Review and Future Research Agenda," Journal of Healthcare Management, vol. 66, no. 2, pp. 150-160, 2021.

[90] R. Kumar, S. Singh, and J. Kaur, "Systematic Analysis of Healthcare Big Data Analytics for Efficient Care and Disease Diagnosing," Health Information Science and Systems, vol. 9, Article 25, 2021. Available: https://doi.org/10.1007/s13755-021-00145-1.

[91] R. K. Ayyasamy, B. Tahayna, S. Alhashmi, S. Eu-Gene, and S. Egerton, "Mining Wikipedia Knowledge to Improve Document Indexing and Classification," in Proc. 10th International Conference on Information Science, Signal Processing and their Applications (ISSPA 2010), 2010.

[92] P. Chinnasamy, D. Roja Ramani, R. K. Ayyasamy, B. J. A. Jebamani, S. Dhanasekaran, and V. Praveena, "Applications of Blockchain Technology in Modern Education System – Systematic Review," in Proc. 2023 International Conference on Computer Communication and Informatics (ICCCI), Jan. 23-25, 2023, Coimbatore, India.

[93] M. Rehman and B. S. Khan, "Automatic Sentiment Annotation of Idiomatic Expressions for Sentiment Analysis Task," Journal of Computer Science and Technology, vol. 28, no. 3, pp. 509-521, May 2013.

[94] L. Smith and J. Doe, "Lexicon-based Non-Compositional Multiword Augmentation: Enriching Tweet Sentiment Analysis," presented at the IEEE International Conference on Data Mining Workshops (ICDMW), Beijing, 2022.

[95] S. Johnson, R. Kumar, and A. Smith, "The Essentials of Sentiment Analysis and Opinion Mining in Social Media: Introduction and Survey of Recent Approaches and Techniques," Journal of Web Science, vol. 6, no. 2, pp. 137-154, Feb. 2021.

[96] C. Lee, J. Kim, and R. K. Gupta, "Applications and Challenges of Implementing Artificial Intelligence in Medical Diagnostics: A Comprehensive Review," Medical Science Monitor, vol. 26, e924957, Nov. 2020.

# REFERENCES

[97] D. Green, L. M. Black, and H. S. Lee, "Enhancing Security Measures in IoT Devices Using Blockchain Technology," in Proc. IEEE Symposium on Security and Privacy Workshops, San Francisco, CA, USA, 2021.

[98] A. Brown and E. Thomas, "Novel Approaches to Data Management and Analytics in Large-Scale Organizations," Journal of Big Data, vol. 8, no. 1, pp. 22-39, Jan. 2021.

Bachelor of Information Systems (Honours) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

# FINAL YEAR PROJECT WEEKLY REPORT
*(Project II)*

| Trimester, Year: Trimester 2 Year 5 | Study week no.: 2 |
|---|---|
| Student Name & ID: DAVID TAN CHOW MENG | |
| Supervisor: DR. RAMESH KUMAR AYYSAMY | |
| Project Title: Challenges and Opportunities in Big Data Analytics: Such as the risks and pitfalls of ignoring context/contextualisation | |

**1. WORK DONE**

Drafting and planning on what to do the next step by continuing from Project 1

**2. WORK TO BE DONE**

Read 50 more journal articles to improve my literature review and make necessary adjustment in chapter 3

**3. PROBLEMS ENCOUNTERED**

No problem encountered.

**4. SELF EVALUATION OF THE PROGRESS**

Put more effort in reading more journal articles.

*ramesh*

_____

Supervisor's signature

_____

Student's signature

A-1

# FINAL YEAR PROJECT WEEKLY REPORT
*(Project II)*

| Trimester, Year: Trimester 2 Year 5 | Study week no.: 4 |
|---|---|
| Student Name & ID: DAVID TAN CHOW MENG | |
| Supervisor: DR. RAMESH KUMAR AYYSAMY | |
| Project Title: Challenges and Opportunities in Big Data Analytics: Such as the risks and pitfalls of ignoring context/contextualisation | |

**1. WORK DONE**

**Done writing my literature review and adjustments in Chapter 3**

**2. WORK TO BE DONE**

**Start experiment phase by studying the datasets and determine the objective of my experiment**

**3. PROBLEMS ENCOUNTERED**

**No problem encountered.**

**4. SELF EVALUATION OF THE PROGRESS**

**Put more effort in understanding the datasets.**

*ramesh*

_____
Supervisor's signature

_____
Student's signature

Bachelor of Information Systems (Honours) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

# FINAL YEAR PROJECT WEEKLY REPORT
*(Project II)*

| Trimester, Year: Trimester 2 Year 5 | Study week no.: 6 |
|---|---|
| **Student Name & ID: DAVID TAN CHOW MENG** | |
| **Supervisor: DR. RAMESH KUMAR AYYSAMY** | |
| **Project Title: Challenges and Opportunities in Big Data Analytics: Such as the risks and pitfalls of ignoring context/contextualisation** | |

**1. WORK DONE**

Studied and understand the datasets and determine the objective.

**2. WORK TO BE DONE**

Start the first half of experiment phase by following the methodology such as Data cleaning, Data integration, visualize data, correlation analysis and model selection

**3. PROBLEMS ENCOUNTERED**

No problem encountered.

**4. SELF EVALUATION OF THE PROGRESS**

Put more effort in my experiment

*ramesh*

_____
Supervisor's signature

_____
Student's signature

A-3

Bachelor of Information Systems (Honours) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

# FINAL YEAR PROJECT WEEKLY REPORT
*(Project II)*

| Trimester, Year: Trimester 2 Year 5 | Study week no.: 8 |
|---|---|
| Student Name & ID: DAVID TAN CHOW MENG | |
| Supervisor: DR. RAMESH KUMAR AYYSAMY | |
| Project Title: Challenges and Opportunities in Big Data Analytics: Such as the risks and pitfalls of ignoring context/contextualisation | |

**1. WORK DONE**
[Please write the details of the work done in the last fortnight.]

**Done with first half of experiment phase by following the methodology such as Data cleaning, Data integration, visualize data, correlation analysis and model selection.**

**2. WORK TO BE DONE**

**Start the second half of experiment which is Model training and tuning, model evaluation, and visualizing and reporting.**

**3. PROBLEMS ENCOUNTERED**

**No problem encountered.**

**4. SELF EVALUATION OF THE PROGRESS**

**Put more effort in my experiment.**

*ramesh*

_____
Supervisor's signature

_____
Student's signature

Bachelor of Information Systems (Honours) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

# FINAL YEAR PROJECT WEEKLY REPORT
*(Project II)*

| Trimester, Year: Trimester 2 Year 5 | Study week no.: 10 |
|---|---|
| **Student Name & ID: DAVID TAN CHOW MENG** | |
| **Supervisor: DR. RAMESH KUMAR AYYSAMY** | |
| **Project Title: Challenges and Opportunities in Big Data Analytics: Such as the risks and pitfalls of ignoring context/contextualisation** | |

**1. WORK DONE**

**Done with experiment**

**2. WORK TO BE DONE**

**Start writing my report.**

**3. PROBLEMS ENCOUNTERED**

**No problem encountered.**

**4. SELF EVALUATION OF THE PROGRESS**

**Put more effort in writing report**

*ramesh*

_____ _____          _____

Supervisor's signature          Student's signature

Bachelor of Information Systems (Honours) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

# FINAL YEAR PROJECT WEEKLY REPORT
*(Project II)*

| Trimester, Year: Trimester 2 Year 5 | Study week no.: 12 |
|---|---|
| **Student Name & ID: DAVID TAN CHOW MENG** ||
| **Supervisor: DR. RAMESH KUMAR AYYSAMY** ||
| **Project Title: Challenges and Opportunities in Big Data Analytics: Such as the risks and pitfalls of ignoring context/contextualisation** ||

**1. WORK DONE**

**Done with all the report writing and amendments.**

**2. WORK TO BE DONE**

**Combine my report into the template given by the university.**

**Send the report to my supervisor for final checking before submission on 26 April 2024**

**3. PROBLEMS ENCOUNTERED**

**No problem encountered.**

**4. SELF EVALUATION OF THE PROGRESS**

**I've did my best on the report**

*ramesh*

_____          _____
Supervisor's signature                      Student's signature

Bachelor of Information Systems (Honours) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

**POSTER**



# Challenges and Opportunities in Big Data Analytics
## The Risks and Pitfalls of Ignoring Context
By
David Tan Chow Meng

**UTAR**
UNIVERSITI TUNKU ABDUL RAHMAN

## Introduction

**1**
- Big Data Analytics plays a transformative role in decision-making across various sectors
- Ignoring context in big data can lead to misleading insights and poor decisions

## Research Objectives

**2**
- To identify and discuss the challenges in big data analytics due to ignored contextual factors
- To highlight the opportunities arising from proper contextualization in big data

## Methodology

**3**
- Use of a dataset from Kaggle for experimental analysis
- Methodological approach: Data collection, preprocessing, exploration, visualization, and model evaluation

## Key Insights

**4**
- Contextualization enhances decision-making effectiveness and data relevance
- Context-aware systems significantly improve user experience and operational efficiency

## Significant Case Examples

**5**
- Application of context-aware systems in smart cities, healthcare, and network orchestration

## Recommendation

**6**
- Integration of pervasive computing and knowledge management to enhance context awareness in big data systems
- Development of innovative algorithms tailored to specific contexts to optimize decision-making processes

## Visuals

**7**
- Graphs showing the impact of context-aware systems on decision-making accuracy
- Diagram of context-aware recommender system algorithm

## Conclusion

**8**
- Emphasizes the importance of context in unblocking the full potential of big data analytics
- Calls for a shift towards more sophisticated, context-aware analytical frameworks

Bachelor of Information Systems (Honours) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR

# Plagiarism Report

CHALLENGES AND OPPORTUNITIES IN BIG DATA
ANALYTICS_SUCH AS THE RISKS AND PITFALLS OF IGNORING
CONTEXT CONTEXTUALISATION.docx

ORIGINALITY REPORT

| 3% | 2% | 2% | 0% |
|---|---|---|---|
| SIMILARITY INDEX | INTERNET SOURCES | PUBLICATIONS | STUDENT PAPERS |

PRIMARY SOURCES

| 1 | hcai.ca.gov<br>Internet Source | 1% |
|---|---|---|
| 2 | "Intelligent Systems Design and Applications",<br>Springer Science and Business Media LLC,<br>2020<br>Publication | <1% |
| 3 | es.scribd.com<br>Internet Source | <1% |
| 4 | smartyads.com<br>Internet Source | <1% |
| 5 | ActEd<br>Publication | <1% |
| 6 | "Intelligent Computing and Networking",<br>Springer Science and Business Media LLC,<br>2023<br>Publication | <1% |
| 7 | Khattak, Asad, Noman Akbar, Mohammad<br>Aazam, Taqdir Ali, Adil Khan, Seokhee Jeon, | <1% |

Myunggwon Hwang, and Sungyoung Lee. "Context Representation and Fusion: Advancements and Opportunities", Sensors, 2014.
Publication

| 8 | Submitted to Swinburne University of Technology<br>Student Paper | <1% |
| 9 | etheses.whiterose.ac.uk<br>Internet Source | <1% |
| 10 | dif7uuh3zqcps.cloudfront.net<br>Internet Source | <1% |
| 11 | "Machine Learning-Based Short-Term Prediction of Air-Conditioning Load through Smart Meter Analytics", Energies, 2017<br>Publication | <1% |
| 12 | Submitted to Addis Ababa University<br>Student Paper | <1% |
| 13 | ebin.pub<br>Internet Source | <1% |
| 14 | erepository.uonbi.ac.ke<br>Internet Source | <1% |
| 15 | eprints.utar.edu.my<br>Internet Source | <1% |
| 16 | www.mdpi.com<br>Internet Source | <1% |

| | | |
|---|---|---|
| 17 | portail-qualite.public.lu<br>Internet Source | <1% |
| 18 | Huisi Wu, Zebin Zhao, Jiafu Zhong, Wei Wang, Zhenkun Wen, Jing Qin. "PolypSeg+: A Lightweight Context-Aware Network for Real-Time Polyp Segmentation", IEEE Transactions on Cybernetics, 2022<br>Publication | <1% |
| 19 | Lecture Notes in Computer Science, 2005.<br>Publication | <1% |
| 20 | Submitted to Universiti Malaysia Pahang<br>Student Paper | <1% |
| 21 | Zhe Qu, Rui Duan, Lixing Chen, Jie Xu, Zhuo Lu, Yao Liu. "Context-Aware Online Client Selection for Hierarchical Federated Learning", IEEE Transactions on Parallel and Distributed Systems, 2022<br>Publication | <1% |
| 22 | dev.to<br>Internet Source | <1% |
| 23 | osuva.uwasa.fi<br>Internet Source | <1% |
| 24 | web.archive.org<br>Internet Source | <1% |
| 25 | www.coursehero.com<br>Internet Source | <1% |

| 26 | Lecture Notes in Computer Science, 2006.<br>Publication | <1 % |
|----|--------------------------------------------------------|------|
| 27 | hanheng Li, Wenyi Zhao, huihua yang.<br>"Enhanced Autofocusing with a Compact and<br>Swift ST-VGG Network", Optica Publishing<br>Group, 2024<br>Publication | <1 % |

Exclude quotes          On                    Exclude matches          Off
Exclude bibliography    Off

PLAGIARISM CHECK RESULT

| Form Title: Supervisor's Comments on Originality Report Generated by Turnitin for Submission of Final Year Project Report (for Undergraduate Programmes) | | | |
|---|---|---|---|
| Form Number: FM-IAD-005 | Rev No.: 0 | Effective Date: 01/10/2013 | Page No.: 1of 1 |

**UTAR**

## FACULTY OF INFORMATION AND COMMUNICATION TECHNOLOGY

| Full Name(s) of Candidate(s) | DAVID TAN CHOW MENG |
|---|---|
| ID Number(s) | 19ACB06731 |
| Programme / Course | INFORMATION SYSTEM ENGINEERING |
| Title of Final Year Project | Challenges and Opportunities in Big data analytics: Such as the risks and pitfalls of ignoring context/contextualisation. |

| **Similarity** | **Supervisor's Comments** **(Compulsory if parameters of originality exceed the limits approved by UTAR)** |
|---|---|
| **Overall similarity index:** ___3___ % **Similarity by source** Internet Sources: ___2___ % Publications: ___2___ % Student Papers: ___0___ % | |
| **Number of individual sources listed** of more than 3% similarity: _0_____ | |
| **Parameters of originality required, and limits approved by UTAR are as Follows:** (i) Overall similarity index is 20% and below, and (ii) Matching of individual sources listed must be less than 3% each, and (iii) Matching texts in continuous block must not exceed 8 words *Note: Parameters (i) – (ii) shall exclude quotes, bibliography and text matches which are less than 8 words.* | |

Note: Supervisor/Candidate(s) is/are required to provide softcopy of full set of the originality report to Faculty/Institute

*Based on the above results, I hereby declare that I am satisfied with the originality of the Final Year Project Report submitted by my student(s) as named above.*

*ramesh*

_____          _____
  Signature of Supervisor                           Signature of Co-Supervisor

  Name: ___Dr. Ramesh Kumar                Name: _____
Ayyasamy___

  Date: 26 April 2024                              Date: _____

Bachelor of Information Systems (Honours) Information Systems Engineering
Faculty of Information and Communication Technology (Kampar Campus), UTAR