PRESERVING CULTURAL HERITAGE SITES THROUGH RANDOM FOREST
AND XGBOOST ALGORITHM FOR MICROCLIMATE MONITORING AND
PREDICTION

IMAN AIDI ELHAM BIN HAIRUL NIZAM

UNIVERSITI TEKNOLOGI MALAYSIA

## UNIVERSITI TEKNOLOGI MALAYSIA

---

### DECLARATION OF THESIS / UNDERGRADUATE PROJECT REPORT AND COPYRIGHT

Author's full name   : Iman Aidi Elham bin Hairul Nizam

Date of Birth   : 22 January 2000

Title   :  Preserving Cultural Heritage Sites through Random Forest and XGBoost Algorithm for Microclimate Monitoring and Prediction

Academic Session   :

I declare that this thesis is classified as:

| | | |
|---|---|---|
| ☐ | **CONFIDENTIAL** | (Contains confidential information under the Official Secret Act 1972)* |
| ☐ | **RESTRICTED** | (Contains restricted information as specified by the organization where research was done)* |
| ✓ | **OPEN ACCESS** | I agree that my thesis to be published as online open access (full text) |

1. I acknowledged that Universiti Teknologi Malaysia reserves the right as follows:

2. The thesis is the property of Universiti Teknologi Malaysia

3. The Library of Universiti Teknologi Malaysia has the right to make copies for the purpose of research only.

4. The Library has the right to make copies of the thesis for academic exchange.

Certified by:

**SIGNATURE OF STUDENT**

**SIGNATURE OF SUPERVISOR**

A20EC5006

AP. DR. MOHD SHAHIZAN OTHMAN

**MATRIX NUMBER**

**NAME OF SUPERVISOR**

Date:  20 JUNE 2015        Date:   20 JUNE 2015

NOTES : If the thesis is CONFIDENTIAL or RESTRICTED, please attach with the letter from the organization with period and reasons for confidentiality or restriction

"I hereby declare that we have read this thesis and in my opinion this thesis is suffcient in term of scope and quality for the award of the degree of Doctor of Philosophy (Specialization)"


Signature               : _____

Name of Supervisor I   : MOHD SHAHIZAN OTHMAN

Date                  : 25 JUNE 2023

**BAHAGIAN A - Pengesahan Kerjasama***

Adalah disahkan bahawa projek penyelidikan tesis ini telah dilaksanakan melalui kerjasama antara Click or tap here to enter text. dengan Click or tap here to enter text.

Disahkan oleh:

Tandatangan :                                    Tarikh :

Nama :

Jawatan :

(Cop rasmi)

* *Jika penyediaan tesis atau projek melibatkan kerjasama.*

---

**BAHAGIAN B - Untuk Kegunaan Pejabat Sekolah Pengajian Siswazah**

Tesis ini telah diperiksa dan diakui oleh:

Nama dan Alamat Pcmeriksa Luar       **:**




Nama dan Alamat Pcmeriksa Dalam     **:**




Nama Penyelia Lain (jika ada)             **:**




Disahkan oleh Timbalan Pendaftar di SPS:

Tandatangan    :                                    Tarikh :

Nama             :

PRESERVING CULTURAL HERITAGE SITES THROUGH RANDOM FOREST
AND XGBOOST ALGORITHM FOR MICROCLIMATE MONITORING AND
PREDICTION

IMAN AIDI ELHAM BIN HAIRUL NIZAM

A thesis submitted in fulfilment of the
requirements for the award of the degree of
Bachelor of Computer Science (Software Engineering)

School of Computing
Faculty of Engineering
Universiti Teknologi Malaysia

JULY 2023

**DECLARATION**

I declare that this thesis entitled *" Preserving Cultural Heritage Sites through Random Forest and XGBoost Algorithm for Microclimate Monitoring and Prediction"* is the result of my own research except as cited in the references. The thesis has not been accepted for any degree and is not concurrently submitted in candidature of any other degree.

Signature      :   .............................................

Name          :   Iman Aidi Elham bin Hairul Nizam

Date           :   25 JUNE 2023

# DEDICATION

This thesis is dedicated to my parents who taught me to work hard and dream big in life. Thank you to my supervisor Associate Prof Dr Mohd Shahizan Othman for guiding me throughout this thesis. Thank you too to my supportive friends who are also struggling to finish their own thesis and has helped me with this thesis either physically or morally.

# ACKNOWLEDGEMENT

In preparing this thesis, I was in contact with many people, researchers, academicians, and practitioners. They have contributed towards my understanding and thoughts. In particular, I wish to express my sincere appreciation to my main thesis supervisor, Professor Dr. Mohd Shahizan Othman, for encouragement, guidance, critics and friendship. Without his continued support and interest, this thesis would not have been the same as presented here.

I am also indebted to Universiti Teknologi Malaysia (UTM) for funding my Bachelor study. Librarians at UTM also deserve special thanks for their assistance in supplying the relevant literatures.

My fellow undergraduate student should also be recognised for their support. My sincere appreciation also extends to all my colleagues and others who have provided assistance at various occasions. Their views and tips are useful indeed. Unfortunately, it is not possible to list all of them in this limited space. I am grateful to all my family member.

# ABSTRACT

The preservation of heritage architecture is significant in maintaining the cultural heritage of a region. This study focuses on the development of a machine learning-based microclimate monitoring system to assist local authorities in Johor Bahru, Malaysia, with planning preventive maintenance actions for a designated heritage site. The research project comprises obtaining microclimate data, including temperature, humidity, and wind speed, from the Malaysian Meteorological Department (MET Malaysia). In order to optimize the monitoring and prediction process, the performance of two machine learning algorithms, Random Forest and XGBoost, will be compared to determine the most suitable method for microclimate analysis. The project also involves designing and developing an interactive, user-friendly dashboard that displays real-time microclimate data using data visualization tools. The effectiveness of the developed algorithm and dashboard will be tested to evaluate their potential in aiding local authorities with the implementation of more effective maintenance plans for the heritage site. This research aims to contribute to the preservation of cultural heritage sites by utilizing advanced machine learning techniques for microclimate monitoring and prediction, ultimately supporting sustainable and efficient conservation efforts.

# ABSTRAK

Pemeliharaan senibina warisan memegang peranan penting dalam mengekalkan warisan budaya suatu kawasan. Kajian ini memberi tumpuan kepada pembangunan sistem pemantauan mikroiklim berasaskan pembelajaran mesin untuk membantu pihak berkuasa tempatan di Johor Bahru, Malaysia, merancang tindakan penyelenggaraan pencegahan untuk tapak warisan yang ditetapkan. Projek penyelidikan ini merangkumi mendapatkan data mikroiklim, termasuk suhu, kelembapan, dan kelajuan angin, daripada Jabatan Meteorologi Malaysia (MET Malaysia). Untuk mengoptimumkan proses pemantauan dan ramalan, prestasi dua algoritma pembelajaran mesin, Random Forest dan XGBoost, akan dibandingkan untuk menentukan kaedah yang paling sesuai untuk analisis mikroiklim. Projek ini juga melibatkan reka bentuk dan pembangunan papan pemuka interaktif dan mesra pengguna yang memaparkan data mikroiklim masa nyata menggunakan alat visualisasi data. Keberkesanan algoritma dan papan pemuka yang dibangunkan akan diuji untuk menilai potensi mereka dalam membantu pihak berkuasa tempatan melaksanakan rancangan penyelenggaraan yang lebih berkesan untuk tapak warisan. Penyelidikan ini bertujuan untuk menyumbang kepada pemeliharaan tapak warisan budaya dengan menggunakan teknik pembelajaran mesin yang maju untuk pemantauan dan ramalan mikroiklim, yang pada akhirnya menyokong usaha-usaha konservasi yang mampan dan cekap.

# TABLE OF CONTENTS

| | TITLE | PAGE |
|---|---|---|

# LIST OF TABLES

# LIST OF FIGURES

x

# LIST OF ABBREVIATIONS

ML        -        Machine Learning

AI        -        Artificial Intelligence

RF        -        Random Forest

CH        -        Cultural Heritage

CSV       -        Comma-separated Values

LR        -        Logistic Regression

ANN       -        Artificial Neural Network

CNN       -        Convolutional Neural Network

KNN       -        K-Nearest Neighbour

# CHAPTER 1

## INTRODUCTION

### 1.1    Overview

Cultural heritage sites are the basis for our global and historical values. They connect us to the traditions left by our ancestors and contribute significantly to the cultural identity of human society (Lombardo et al., 2020). The preservation of cultural heritage, whether it be buildings or artifacts, is subject to various risks of damage and deterioration that result from microclimate conditions in the surrounding environment. These conditions are determined by several factors, including microclimate parameters such as temperature, humidity, airborne pollutants concentrations, air speed, and others (Fabbri & Bonora, 2021). Particularly in developing nations, these impacts pose a significant challenge to the preservation of cultural heritage (Pioppi et al., 2020). Safeguarding worldwide cultural heritage sites is of utmost importance for preserving cultural identity and human heritage, as well as promoting cultural and tourism-driven economic development (Alcaraz Tarragüel et al., 2012).

In recent years, the administration of cultural heritage sites and monuments has gained worldwide focus through the implementation of detection, monitoring, and comprehensive assessment methods. Initiatives are also underway to enhance and preserve these heritage resources by adopting suitable adaptation measures and sustainable management approaches (Guzman et al., 2020). To address these challenges, this thesis focuses on the application of advanced machine learning algorithms, namely Random Forest and XGBoost, for microclimate monitoring and prediction at cultural heritage sites. By leveraging these techniques, it aims to contribute to the preservation of cultural heritage sites under changing environmental conditions, ultimately supporting sustainable and efficient conservation efforts.

1

## 1.2 Problem Background

Cultural heritage sites have consistently drawn visitors who seek to spend quality time and pursue unique experiences by engaging with local cultures and communities (Ramkissoon et al., 2013). As a result, the economies of these tourist destinations largely rely on attracting visitors, encouraging repeat visits, garnering recommendations, and generating positive word-of-mouth regarding the locations (Rezapouraghdam et al., 2021). In addition, the natural environments in which tourism activities occur are also enhancing the well-being and quality of life for local residents (Ramkissoon et al., 2018). Lately, Johor Bahru has been experiencing frequent climate fluctuations that negatively impact the aesthetic appeal of the area's heritage sites, significantly affecting the industry of tourism and local economy. Generally, microclimate changes in these regions cause substantial damage to cultural heritage sites and various monuments. Consequently, striking a balance between consumption and conservation strategies presents increasing challenges for the effective management of cultural heritage sites (Buonincontri et al., 2017). Therefore, focusing on the preservation of cultural heritage and promoting sustainable tourism has become a primary objective recently to support both cultural heritage tourism and the overall well-being of communities (Megeirhi et al., 2020).

## 1.3 Research Aim

The goal of this study is to analyze vulnerable zones of cultural heritage (CH) sites and monuments in Johor Bahru, Malaysia, by employing microclimate monitoring and prediction through the Random Forest and XGBoost algorithms. By assessing temperature, humidity, and wind speed, the study aims to maintain environmental sustainability at these heritage sites. In this research, we have prepared a microclimate monitoring dashboard and evaluated the significance of factors contributing to microclimate changes. The Random Forest and XGBoost algorithms were employed to analyze the impact of these factors on the preservation of CH sites.

## 1.4    Research Objectives

The following are the objectives proposed:

(a)    To investigate and identify the most suitable machine learning algorithms for analyzing microclimate data, and to recognize patterns, trends, and potential issues related to the heritage site's preservation.

(b)    To evaluate the effectiveness of the developed machine learning models and the dashboard in assisting local authorities to plan preventive maintenance actions that preserve the site's aesthetics and cultural values.

(c)    To develop and design a user-friendly dashboard that displays real-time microclimate data and provides recommendations for maintenance actions to assist local authorities with heritage site preservation.

## 1.5    Research Scopes

The scope of this research project covers several aspects related to the preservation of the Johor Bahru High Court and Sultan Ibrahim Building in Johor Bahru, using machine learning-based microclimate monitoring. The primary focus is on the development of a dashboard to collect, display, and analyze microclimate parameters for assisting the local authority in planning preventive maintenance actions for these two heritage sites. Specific areas included in the scope of this research are:

(a)    Research involves obtaining microclimate data from the Malaysian Meteorological Department (MET Malaysia) for a designated heritage site in Johor Bahru. This data contains parameters like temperature, humidity, and wind speed.

(b)    The research intends to compare the performance of two different machine learning algorithms between Random Forest and XGBoost to determine the most suitable method for microclimate monitoring and prediction.

(c)     The project includes designing and developing an interactive user-friendly dashboard by using data visualization tools that display real-time microclimate data.

(d)     The research will involve testing the effectiveness of the developed algorithm and dashboard in assisting local authorities with planning more effective maintenance plans for the heritage site.

**1.6     Research Contribution**

A thorough literature review on microclimate impacts on CH sites reveals that many researchers have used various statistical and machine learning methods, including Logistic Regression (LR), Artificial Neural Network (ANN), Convolutional Neural Network (CNN), K-Nearest Neighbor (KNN), and Support Vector Machine (SVM), to create microclimate monitoring and prediction dashboards. However, the combination of Random Forest and XGBoost algorithms, along with the analysis of temperature, humidity, and wind speed, has not yet been employed in the context of CH site preservation. As a result, this study offers a novel contribution to the machine learning field, particularly for modeling microclimate threats and risk assessments of cultural heritage sites.

In the context of the current changing climate and landscape, this study is highly relevant and makes a significant contribution to sustainable management of cultural heritage resources. Climate change can pose a significant threat to the integrity of heritage sites due to its impact on key environmental factors such as temperature, humidity, and wind patterns. This can lead to increased vulnerability and potential damage to these cultural assets. The study offers valuable insights and technical guidance regarding the selection of input causative factors, appropriate machine learning algorithms, and proper interpretation and evaluation of outcomes, which can inform future research and decision-making processes.

Moreover, this study has essential implications for the conservation of natural resources and heritage sites in Johor Bahru, Malaysia. The findings of this study are expected to have practical applications for professionals involved in land use planning,

landscape management, archaeological preservation, and public administration, as they strive to effectively manage cultural heritage sites and promote environmental sustainability through evidence-based strategies. By monitoring and predicting microclimate changes using Random Forest and XGBoost algorithms, stakeholders can better preserve and protect cultural heritage sites for future generations.

## 1.7    Report Organization

This report comprises five chapters. Chapter 1 introduces the topic of preserving cultural heritage sites through microclimate monitoring and prediction using Random Forest and XGBoost algorithms, the research background, and the purpose of conducting this study in Johor Bahru, Malaysia. Chapter 2 discusses the literature review related to microclimate monitoring, the assessment of temperature, humidity, and wind speed, as well as the comparison of machine learning techniques for processing and analyzing data from heritage sites. Chapter 3 delves into methodology of the research, describing how the study is conducted using the Random Forest and XGBoost algorithms to measure and analyze the data on temperature, humidity, and wind speed for preserving cultural heritage sites. Chapter 4 presents the research design and implementation, detailing how the experiment was executed to extract valuable insights from the microclimate data. Chapter 5 showcases the results obtained from this research, including the Power BI dashboard displaying the analyzed data. Finally, Chapter 6 offers a summary and conclusion for the study, highlighting the key findings and implications for the preservation of cultural heritage sites through microclimate monitoring and prediction using Random Forest and XGBoost algorithms.

# CHAPTER 2

# LITERATURE REVIEW

## 2.1    Introduction to Case Study

Climate change has emerged as a significant global challenge in recent years, impacting various sectors, including the preservation of cultural heritage sites. The increasing frequency and intensity of extreme weather events, along with gradual shifts in temperature, humidity, and wind patterns, have highlighted the need for adaptive solutions to safeguard these invaluable assets.

One promising approach to address these challenges involves the application of advanced algorithms, such as Random Forest and XGBoost, for microclimate monitoring and prediction at cultural heritage sites. These techniques can help preserve and protect these valuable assets by analysing temperature, humidity, and wind speed data, which are crucial factors in the conservation of these sites.

This case study focuses on the implementation of Random Forest and XGBoost algorithms for microclimate monitoring and prediction at two cultural heritage sites in Johor Bahru, Malaysia: the Johor Bahru High Court and the Sultan Ibrahim Building. By leveraging these advanced techniques and utilizing Power BI dashboards for data visualization and analysis, this research aims to enhance the understanding of site-specific microclimates and inform effective conservation strategies for these historic landmarks.

Through continuous assessment and refinement of these methods, researchers, conservators, and heritage site managers can work together to develop improved strategies for preserving cultural heritage sites like the Johor Bahru High Court and the Sultan Ibrahim Building under changing environmental conditions. By adopting a collaborative approach, we can ensure the protection and preservation of these

invaluable assets for future generations to appreciate and learn from, even in the face of challenges posed by climate change.

## 2.2    Importance of Preserving Cultural Heritage Sites

The preservation of cultural heritage sites holds immense significance for society, history, and identity, as these sites serve as tangible reminders of our shared past, providing valuable insights into the cultural, social, and economic development of human civilizations (Lowenthal, 1985). By protecting and maintaining these sites, we ensure the continuity of our cultural memory and allow future generations to appreciate and learn from the rich tapestry of human history (UNESCO, 1972). Moreover, cultural heritage sites contribute to a sense of belonging and pride within communities, fostering social cohesion and promoting intercultural dialogue (Smith, 2006). Furthermore, preserving these sites can offer economic benefits, as they often attract tourism and stimulate local economies (Timothy & Boyd, 2003). Given these multifaceted advantages, it is crucial to develop and implement strategies to safeguard cultural heritage sites against various threats, including the impact of microclimate factors, to ensure their longevity and continued cultural relevance.

## 2.3    Impact of Microclimate Factors on Cultural Heritage Sites

Microclimate factors, such as humidity, temperature, wind speed, play a significant role in the deterioration of cultural heritage sites. Existing studies have established the adverse effects of these factors on various materials and structures, leading to both physical and chemical degradation (Cassar, 2005; Camuffo, 2014).

Temperature fluctuations, especially in the presence of moisture, can lead to the expansion and contraction of materials like stone, brick, and mortar, resulting in cracks, delamination, and structural damage (Camuffo, 2014). Moreover, extreme temperatures can accelerate the decay of organic materials, such as wood and textiles, commonly found in cultural heritage sites (Cassar, 2005).

8

Humidity is another critical factor in the deterioration process. High humidity levels can cause moisture to accumulate in porous materials, leading to the growth of mold and bacteria, which can weaken and damage the structure (Lankester & Brimblecombe, 2012). Additionally, the presence of moisture can facilitate the dissolution of soluble salts in porous materials, causing efflorescence and sub florescence, further compromising structural integrity (Cassar, 2005).

Wind speed, particularly in combination with rain, can exacerbate the erosion of building materials and increase the rate of material loss from structures (Cassar, 2005). Moreover, high wind speeds can cause physical damage to fragile elements, such as decorative features and stained-glass windows (Camuffo, 2014).

In summary, understanding the impact of microclimate factors on cultural heritage sites is crucial for developing effective preservation strategies. By identifying and mitigating the risks associated with these factors, we can better protect these invaluable resources and ensure their continued existence for future generations.

## 2.4    Traditional Methods for Cultural Heritage Sites Preservation

Traditional methods for cultural heritage site preservation often rely on reactive maintenance approaches. These approaches involve responding to issues and damage after they have already occurred, rather than anticipating and preventing them. Reactive maintenance has several limitations, making it necessary to explore proactive and preventive measures for the preservation of cultural heritage sites (Staniforth, 2013).

Delayed intervention is one of the limitations of reactive maintenance, as it occurs after the damage has been detected, leading to further deterioration or irreversible loss of cultural elements (Muñoz Viñas, 2002). Additionally, reactive maintenance can be expensive, especially if the damage requires extensive interventions and specialized expertise (Stovel, 2005). Incomplete recovery can also be an issue, as advanced damage can result in the loss of original features or materials, compromising the site's authenticity and historical value (Muñoz Viñas, 2002).

Moreover, interventions during reactive maintenance can be invasive or destructive, leading to further damage or exposing other areas to new risks (Matero, 1999).

To address these limitations, there is a need to shift towards proactive and preventive measures, such as regular monitoring, preventive conservation, maintenance planning, and capacity building for local stakeholders. Regular inspections and monitoring can help identify early signs of deterioration or potential threats (Caple, 2008), while preventive conservation can reduce or eliminate risk factors contributing to the site's deterioration, such as controlling humidity and temperature (Muñoz Viñas, 2002). Maintenance planning, including preventive measures and timely interventions, can also help address potential issues (Caple, 2008). Capacity building through training and education for local stakeholders can further enhance the site's preservation efforts (Ashley-Smith, 2016).

In conclusion, the preservation of cultural heritage sites requires a shift towards proactive and preventive measures, which can minimize the risk of irreversible damage, maintain the site's authenticity and historical value, and reduce the overall cost of preservation efforts. By adopting regular monitoring, preventive conservation, maintenance planning, and capacity building for local stakeholders, cultural heritage sites can be better preserved for future generations (Muñoz Viñas, 2002; Caple, 2008; Ashley-Smith, 2016).

## 2.5    Machine Learning Algorithms

This study develops machine learning-based methods for microclimate monitoring and prediction at cultural heritage sites, using the supervised learning concept. This involves training a classifier to assign labels to specific data points or regions in the dataset, enabling it to identify hidden patterns and signatures of various labelled factors and make accurate predictions. To ensure effective monitoring and prediction using a variety of data sources, it is crucial to use classifiers that can handle large-scale data and achieve high accuracy quickly. The study focuses on two classifiers, XGBoost and Random Forest, which are both capable of achieving these requirements.

### 2.5.1 Random Forest



**Figure 1: Random Forest Model Architecture**

Breiman's Random Forest algorithm, introduced in 2001, is a widely used ensemble learning model that is known for its versatility in performing various tasks such as classification, regression, clustering, interaction detection, and variable selection (Rahmati et al., 2017; Belgiu and Drăguţ, 2016). This learning method leverages the aggregation of decision trees, which divide input data based on specific parameters in a tree-like structure (Ma and Cheng, 2016; Breiman, 2001) (see Fig. 1). Unlike other learning methods, Random Forest is designed to handle complex datasets with high dimensionality, noisy, and missing data, making it particularly useful for microclimate monitoring and prediction at cultural heritage sites.

Each decision tree in a Random Forest model is built using a bootstrapped sample of the data, with nodes split according to the optimal subset and randomly selected predictors at each stage (Araki et al., 2018; Rahmati et al., 2017). The final classification is based on the majority vote of the decision trees, and output is generated accordingly (Micheletti et al., 2014; Rahmati et al., 2017). This approach helps prevent overfitting, where a model learns the training data too well and fails to generalize well to new data. Random Forest's robustness and high-performance capabilities have made

it a popular choice in various fields, including image analysis, remote sensing, and ecology.

Furthermore, Random Forest is a highly flexible model that can handle a wide range of input variables, including categorical and continuous variables, and can deal with missing data points by imputing values. The algorithm can also determine the importance of each input variable in predicting the output, enabling researchers to identify the most influential parameters for microclimate monitoring and prediction at cultural heritage sites. Additionally, researchers have developed various extensions and modifications to improve the algorithm's efficiency, such as parallel computing, pruning techniques, and feature importance measures (Balogun et al., 2021; Tella et al., 2021).

One notable feature of Random Forest is its ability to handle interactions between variables, which is important in predicting microclimate parameters at cultural heritage sites. The algorithm can identify and model complex interactions between multiple variables, allowing researchers to better understand the relationships between environmental factors and microclimate patterns. This feature is particularly valuable in cultural heritage sites, where environmental conditions can vary significantly and interactions between environmental factors can be complex.

In conclusion, Random Forest is a powerful machine learning algorithm that has proven to be a valuable tool for microclimate monitoring and prediction at cultural heritage sites. Its robustness, flexibility, and high-performance capabilities make it an attractive choice for handling complex datasets with high dimensionality, noisy, and missing data. Moreover, the algorithm's ability to identify and model interactions between variables provides researchers with valuable insights into the complex relationships between environmental factors and microclimate patterns.

### 2.5.2 XGBoost



**Figure 2: XGBoost Model Architecture**

XGBoost is a popular machine learning algorithm that is commonly used for classification tasks. It belongs to the family of boosting algorithms, where multiple weak learners are combined to create a strong model. The algorithm works by iteratively adding decision trees to the model and adjusting their weights based on the error rate of the previous trees. The result is a highly accurate classifier that can handle large and complex datasets.

One of the key advantages of XGBoost is its ability to handle missing data effectively. The algorithm can use surrogate splits to compensate for missing data points, resulting in improved accuracy and robustness in the presence of missing data. XGBoost is also highly optimized for parallel computing, enabling it to process large volumes of data quickly and efficiently.

XGBoost has demonstrated high performance and accuracy when dealing with large-scale, multi-class data in various fields, including remote sensing, medical diagnosis, and natural language processing. Studies have shown that XGBoost can outperform other popular classification algorithms, such as Random Forest and

Support Vector Machines (SVM), in terms of accuracy and efficiency (Bhagwat & Shankar, 2019; Zamani Joharestani et al., 2019; Rumora et al., 2020). This makes it an attractive choice for microclimate monitoring and prediction at cultural heritage sites, where large datasets and high-dimensional feature spaces are common.

XGBoost is highly scalable, which makes it an ideal choice for handling large volumes of satellite data. This enables researchers to perform microclimate monitoring and prediction in real-time, providing valuable insights into the environmental conditions at cultural heritage sites. Additionally, XGBoost is highly optimized for feature selection, allowing researchers to identify the most influential variables for microclimate monitoring and prediction.

In conclusion, XGBoost is a powerful machine learning algorithm that offers several unique advantages for microclimate monitoring and prediction at cultural heritage sites. Its ability to handle missing data, parallel computing, scalability, and feature selection capabilities make it an attractive choice for researchers and practitioners in this field. By leveraging XGBoost's powerful capabilities, researchers can gain valuable insights into the environmental conditions at cultural heritage sites, enabling them to develop more effective strategies for managing and preserving these invaluable assets for future generations.

## 2.6    Comparative Analysis of Previous Case Studies and the Uses of Machine Learning in Cultural Heritage Preservation

Several studies in the cultural heritage field apply machine learning (ML) techniques for tasks such as automatic text recognition, image annotation, and user preference recommendations. However, the use of ML in conservation science and heritage preservation studies is limited. These studies primarily focus on identifying and classifying materials or structures or using ML to monitor cultural heritage collections or sites for abnormalities. For instance, Zou et al. employed deep learning on image data to locate missing or damaged heritage components in historical buildings, while Kejser et al. used ML to classify the acidity of historic paper samples. Pei et al. utilized machine learning to predict household mite infestation based on

indoor climate conditions and found that the extreme gradient boosting (XGBoost) model was the most suitable approach.

Table 1: Comparative Analysis of Previous Case Studies and the Uses of Machine Learning in Cultural Heritage Preservation

| Case Study | Method Used | Target Site/Subject | Main Outcomes |
|---|---|---|---|
| Yu et al. (2022) | Convolutional Neural Network Deep Learning | Dunhuang Mogao Grottoes, China | Detected wall painting deterioration; informed preventive measures |
| (Kumar et al. (2019) | Logistic Regression, Support Vector Machine | Damaged Heritage Sites from 2015 Nepal Earthquake | Classify heritage and not-heritage sites; damage or no damage |
| Prieto et al., (2017) | Multiple Linear Regression, Fuzzy Logic Models | 100 parish churches, located in Seville, Spain | Identifies relevant variables for the functional degradation of the churches. |
| Gonthier et al., (2019) | Support Vector Machines | Child Jesus, the crucifixion of Jesus, Saint Sebastian | Recognition of iconographic elements in artworks. |
| Valero et al., (2019) | Logistic Regression, Multi Class Classification | Chapel Royal in Stirling Castle, Scotland | Identifies loss of material defects and discoloration on the walls. |

## 2.7    Implementation of Random Forest and XGboost in Microclimate Monitoring and Prediction

Researchers have been exploring the performance of XGBoost and Random Forest classifier algorithms for microclimate monitoring and prediction in various studies. These algorithms have proven to be effective in providing valuable insights for monitoring and managing microclimate factors in different environments.

In a study by J. Angelin Jebamalar & A. Sasi Kumar (2019), a hybrid light tree and light gradient boosting model were used for predicting PM2.5 levels. The proposed method captured PM2.5 data using a sensor with Raspberry Pi and stored it in the cloud, where the hybrid model was used for prediction. The hybrid model outperformed other algorithms, including Linear Regression, Lasso Regression, Support Vector Regression, Neural Network, Random Forest, Decision Tree, and XGBoost. Despite its advantages in handling large amounts of data and requiring less space, the hybrid model's limitation was its time-consuming nature.

In a study by Maryam Aljanabi (2020), the authors compared Multilayer Perceptron, XGBoost, Support Vector Regression, and Decision Tree Regressor to predict ozone levels based on temperature, humidity, wind speed, and wind direction. After pre-processing the data and performing feature selection, XGBoost emerged as the superior model for predicting ozone levels on a day-to-day basis.

Soubhik et al. (2018) compared various algorithms, including Linear Regression, Neural Network Regression, Lasso Regression, ElasticNet Regression, Decision Forest, Extra Trees, Boosted Decision Tree, XGBoost, K-Nearest Neighbor, and Ridge Regression, to predict air pollutant levels. They found that XGBoost provided better accuracy due to the arrangement of features in decreasing order of importance for predicting upcoming values. Haotian Jing & Yingchun Wang (2020) used XGBoost to predict the air quality index. By employing weak classifiers and using the shortcomings of previous weak classifiers to form a strong classifier, XGBoost reduced the error between predicted and actual values. However, it was

susceptible to outliers and unwanted air pollutants, as it took the previous value into account.

Mejía et al. (2018) determined PM10 levels best with Random Forest but found that it did not accurately predict the levels of dangerous pollutants. However, Random Forest had the advantage of working with incomplete datasets. Pasupuleti et al. (2020) compared Decision Tree, Linear Regression, and Random Forest for predicting air pollutant levels using meteorological conditions and data from the Arduino platform. Random Forest provided more accurate results due to reduced overfitting and error. However, it required more memory and incurred higher costs.

In summary, XGBoost and Random Forest have been applied in various case studies for microclimate monitoring and prediction, with both algorithms demonstrating their effectiveness in predicting air pollutant levels. While they have their respective limitations, these advanced techniques offer valuable tools for researchers and practitioners seeking to understand and manage the air quality in different environments.

### 2.7.1   Comparison Between Random Forest and XGBoost Algorithms

Table 2: List of Difference between Random Forest and XGBoost Algorithms

| Criteria | XGBoost | Random Forest |
|---|---|---|
| Model Type | Gradient boosting decision tree ensemble (Chen & Guestrin, 2016) | Decision tree ensemble (Breiman, 2001) |
| Learning Approach | Gradient boosting, optimizing loss function (Friedman, 2001) | Bagging, independent decision trees combined through majority voting or averaging (Liaw & Wiener, 2002) |
| Handling Missing Data | Imputation or treating missing values as separate | Imputation or treating missing values as separate categories (Breiman, 2001) |

| | | categories (Chen & Guestrin, 2016) | |
|---|---|---|---|
| Overfitting Prevention | Shrinkage and regularization (Chen & Guestrin, 2016) | Averaging results of multiple decision trees (Breiman, 2001) |
| Interpretability | Can provide feature importance information | Easier to interpret due to simpler decision tree structure (Breiman, 2001) |
| Speed and Scalability | Slower in training due to sequential nature (Chen & Guestrin, 2016) | Faster and more parallelizable due to independent tree construction (Breiman, 2001) |
| Performance | Compare using MAE, RMSE, R-squared (Caruana & Niculescu-Mizil, 2006) | Compare using MAE, RMSE, R-squared (Caruana & Niculescu-Mizil, 2006) |
| Feature Importance | Can rank input variables by importance (Chen & Guestrin, 2016) | Can rank input variables by importance (Breiman, 2001) |
| Hyperparameter Tuning | Requires tuning, may be more sensitive to hyperparameter settings (Probst, Wright, & Boulesteix, 2019) | Requires tuning, may be less sensitive to hyperparameter settings (Probst, Wright, & Boulesteix, 2019) |
| Memory Usage | Less memory usage due to sequential nature (Chen & Guestrin, 2016) | More memory usage due to storage of multiple decision trees (Breiman, 2001) |

## 2.8    Chapter Summary

Through this chapter, the study focuses on the preservation of cultural heritage sites through machine learning-based microclimate monitoring, with a specific focus on the application of Random Forest and XGBoost algorithms at two heritage sites in Johor Bahru, Malaysia. The review begins with an overview of the impact of climate change on cultural heritage sites and the importance of their preservation. It then explores the impact of microclimate factors on cultural heritage sites, including temperature, humidity, and wind speed, and the traditional reactive methods used for

preservation. The limitations of reactive maintenance and the need for a shift towards proactive and preventive measures are discussed, such as regular monitoring and preventive conservation. Finally, the review explains the use of machine learning algorithms in microclimate monitoring and prediction, specifically Random Forest and XGBoost, and their application in this study. The review highlights the significance of using machine learning-based approaches for preserving cultural heritage sites and the potential benefits of incorporating them into preservation strategies.

# CHAPTER 3

# RESEARCH METHODOLOGY

## 3.1    Introduction

This chapter will elaborate more on the methodology used in this research which includes the research design framework, steps and techniques used for the research. The overall research workflow consists of four phases which will be discussed further in this chapter. Each step of the framework will be elaborated in detail on how it will be implemented to this research including the techniques that will be used.

## 3.2    Research Workflow

The research methodology framework consists of four main phases: Literature Review, Data Requirement and Data Collection, Machine Learning Model Development, and Dashboard Development. The framework begins with Phase 1, which involves conducting a literature review on microclimate monitoring and prediction, as well as machine learning techniques (as illustrated in Figure 3). The framework then proceeds to Phase 2, which focuses on identifying the data requirements and collecting the necessary data for addressing the microclimate monitoring and prediction problem. The subsequent phase is the Machine Learning Model Development phase, which entails developing the models. Finally, the framework concludes with the Dashboard Development phase, wherein the obtained results are integrated into a high-fidelity dashboard.

Figure 3: Research Workflow

### 3.2.1 Phase 1: Literature Review

In the first phase, a comprehensive literature review is conducted to gather relevant information on preserving cultural heritage sites through microclimate monitoring and prediction using Random Forest and XGBoost algorithms. Various scholarly sources, including journals, articles, and theses, are explored to understand the current state of research in this field. The literature review delves into topics such as data collection methods, pre-processing techniques, and the application of machine learning models. By examining existing studies, this phase helps identify gaps, challenges, and potential solutions for effectively monitoring and predicting microclimate conditions at heritage sites. The insights gained from the literature review form the foundation for the subsequent phases and guide the research towards developing a robust methodology.

### 3.2.2 Phase 2: Data Requirement and Data Collection

In the second phase, the focus shifts to obtaining microclimate data from the Malaysian Meteorological Department for a specific heritage site in Johor Bahru. This data, encompassing temperature, relative humidity, and wind, rainfall, and solar

radiation measurements, serves as the basis for subsequent analysis and modelling. To ensure data quality, a rigorous pre-processing stage is undertaken, which involves cleaning the raw data and addressing any missing values or outliers. Techniques such as interpolation, statistical analysis, and data imputation are employed to enhance the integrity and accuracy of the collected data. Additionally, feature engineering techniques are applied to extract meaningful features from the raw data, enabling the capturing of temporal dependencies and relationships between variables. This phase prepares the dataset for further analysis and model development in the subsequent phases.

### 3.2.3 Phase 3: Machine Learning Model Development

The third phase revolves around the development of machine learning models for microclimate monitoring and prediction. The pre-processed data is divided into training and testing sets, with the training set used to train and optimize the Random Forest and XGBoost algorithms. Various parameters and hyperparameters are fine-tuned using techniques like grid search and cross-validation to achieve optimal model performance. The trained models are then evaluated using appropriate assessment metrics, such as mean absolute error and mean squared error, to assess their predictive capabilities. This evaluation process helps determine the effectiveness and performance of the Random Forest and XGBoost algorithms in accurately predicting microclimate patterns for the designated heritage site. The models' performance and generalizability are crucial factors in ensuring the reliability and usefulness of the developed models for microclimate monitoring and prediction.

### 3.2.4 Phase 4: Dashboard Development

In the fourth phase, a user-friendly dashboard is designed and developed to visualize and present the microclimate data in a comprehensible manner. The dashboard provides real-time insights into the microclimate conditions of the heritage site, displaying key metrics such as temperature, humidity, and wind speed. The trained machine learning models are integrated into the dashboard to provide recommendations for preventive maintenance actions based on the analyzed data.

Additionally, visualization tools and interactive features are incorporated to facilitate a better understanding of trends and patterns in the microclimate data. The dashboard serves as a valuable tool for local authorities and stakeholders involved in the preservation of cultural heritage sites, enabling them to make informed decisions and take proactive measures to mitigate potential issues that may impact the site's condition.

## 3.3    Justification of Tools, Techniques and Data

The chosen tools for this research include data collection from the Malaysian Meteorological Department (MET Malaysia), data pre-processing techniques, machine learning algorithms (Random Forest and XGBoost), and dashboard development. These tools have been selected based on their suitability for addressing the research objectives and providing valuable insights for preventive maintenance strategies at the designated heritage site in Johor Bahru.

Microclimate data from MET Malaysia is essential for understanding the environmental conditions at the heritage site. Temperature, humidity, and wind speed are crucial parameters that can affect the preservation of heritage structures. By obtaining this data, we can assess the impact of these environmental factors on the site's structural integrity and identify potential maintenance issues.

The research utilizes machine learning techniques to analyze the collected microclimate data and develop predictive models for preventive maintenance. Random Forest and XGBoost algorithms are chosen due to their proven effectiveness in handling complex relationships between variables and handling both regression and classification tasks. These algorithms can capture temporal dependencies in the data and provide accurate predictions for future microclimate conditions.

### 3.3.1   Preservation of Cultural Heritage

The preservation of heritage sites is of paramount importance to maintain cultural identity and historical significance. By utilizing advanced data analysis

techniques, this research aims to provide insights into the impact of microclimate conditions on the designated heritage site in Johor Bahru. The findings can help authorities develop targeted preventive maintenance strategies to ensure the site's long-term preservation.

## 3.4    Data-Driven Decision Making

By collecting and analyzing microclimate data, this research enables data-driven decision making for preventive maintenance. Traditional approaches may not consider the dynamic nature of microclimate conditions and their impact on heritage sites. The use of machine learning algorithms allows for a more comprehensive understanding of the relationships between environmental factors and potential maintenance issues, leading to more informed decision making.

## 3.5    Efficiency and Cost-Effectiveness

The research encourages collaboration between local authorities, heritage site management teams, and relevant stakeholders. By involving these parties in the evaluation and testing phases, their feedback can be gathered to refine the system and ensure its usability and effectiveness. Engaging stakeholders throughout the research process increases their ownership and facilitates the adoption of preventive maintenance strategies.

## 3.6    Real-Time Monitoring and Visualization

The development of a user-friendly dashboard integrating real-time microclimate data and machine learning models provides a powerful tool for monitoring and maintenance planning. The visualization tools within the dashboard help users understand trends and patterns in the data, facilitating proactive decision making. This allows local authorities to respond promptly to changing microclimate conditions and potential threats to the heritage site.

**3.7      Stakeholder Engagement and Collaboration**

The preservation of heritage sites is of paramount importance to maintain cultural identity and historical significance. By utilizing advanced data analysis techniques, this research aims to provide insights into the impact of microclimate conditions on the designated heritage site in Johor Bahru. The findings can help authorities develop targeted preventive maintenance strategies to ensure the site's long-term preservation.

**3.8      Chapter Summary**

This chapter summarized the four phases of the research study. The literature review phase involved a comprehensive review of relevant literature, providing a solid knowledge base for the subsequent phases. The data collection and pre-processing phase focused on obtaining and cleaning microclimate data, while the machine learning model development phase involved training and evaluating Random Forest and XGBoost algorithms. The final phase centred on developing a user-friendly dashboard that visualizes the microclimate data and provides maintenance recommendations.

# CHAPTER 4

# RESEARCH DESIGN AND IMPLEMENTATION

## 4.1    Introduction

This chapter discusses in depth the research design and implementation of the research methodology described in the previous chapter. The proposed solution will be broken down into several steps, including the data collection, pre-processing, feature extraction, training and testing and machine learning models development.

## 4.2    Proposed Solution

The dataset used in this study is collected and extracted from Malaysian Meteorological Department. The dataset that is used in this study mainly from microclimate data such as temperature, humidity, wind speed, rainfall, and solar radiation. The dataset must go through a few experimental steps in order to obtain the sentiment score of the social media data. First, data must be extracted from social media. Following that, the data is subjected to pre-processing and feature extraction. The data then moves on to the next step, which is the feature vector and sentiment classifier implementation step. The final step in this research is to calculate the sentiment score for each tweet. Python code will be used to run each process in this experiment.

**4.3     Research Area Zone Mapping**



Figure 4: Area of Research Zone

In the context of this research, a research area zone mapping was established to focus on the preservation of cultural heritage sites in Johor Bahru, Malaysia. The selected heritage sites for this study are:

1. Sultan Ibrahim Building: This building holds historical and architectural significance as the state secretariat building of Johor. It represents the rich cultural heritage of the region.

2. Johor Bahru High Court: The Johor Bahru High Court is a notable judicial institution that plays a crucial role in the administration of justice. It possesses historical and legal importance.

3. Sultan Abu Bakar Mosque: Known for its exquisite Moorish-inspired architecture, the Sultan Abu Bakar Mosque is a revered religious site. It serves as a symbol of faith and cultural identity.

4. Malayan Railway Museum: The Malayan Railway Museum, located in Johor Bahru, showcases the historical development and significance of the railway system in Malaysia. It offers insights into the country's transportation heritage.

To ensure comprehensive monitoring and analysis of the microclimate conditions in the research area, the chosen auxiliary station is the Sultanah Aminah Hospital. This auxiliary station, situated close to the selected heritage sites, serves as an essential data collection point for microclimate variables. By selecting an auxiliary station in proximity to all the chosen heritage sites, the research study aims to capture accurate and representative microclimate data specific to the designated research area zone.

By focusing on the Sultanah Aminah Hospital auxiliary station, which covers all the selected heritage sites, the research can effectively monitor and analyze the microclimate conditions in the vicinity of these sites. This approach ensures that the data collected and analyzed is directly relevant to the preservation efforts and maintenance strategies of the cultural heritage sites in Johor Bahru.

## 4.4 Experiment Design

### 4.4.1 Microclimate Data Collection Process

In this research study, data collection is a crucial step in developing the prototype dashboard. To ensure the accuracy of the results, data was collected from the Malaysian Meteorological Department (MET Malaysia) website and stored in our database. A duration of 30 years was chosen to gather a substantial amount of data, enabling more reliable and robust predictions. The data was specifically obtained from the auxiliary station of Hospital Johor Bahru, which is in close proximity to the research area.

Five categories of microclimate data were collected for this research, namely temperature, relative humidity, wind, rainfall, and solar radiation. A comprehensive set of attributes was acquired for temperature data, consisting of eight different types.

Similarly, for relative humidity data, five distinct attributes were collected. In the case of rainfall, four attributes were available for analysis. However, due to limited availability of wind data at the Hospital Johor Bahru auxiliary station, only one attribute, namely the hourly surface wind, could be obtained. Consequently, to enhance the predictive capabilities of the research study, the inclusion of solar radiation and rainfall data was deemed necessary. All of the collected data is in the CSV file format and accessible from the MET Malaysia website.

The meticulous collection and inclusion of these various microclimate data categories and attributes aim to ensure a comprehensive analysis and prediction process within the research study. By incorporating multiple data sources, a more holistic understanding of the microclimate conditions at the designated heritage site can be achieved, ultimately facilitating the development of an effective and informative dashboard prototype.

## 4.4.2 Pre-processing and Feature Extraction of Microclimate Data

After the data are collected, they must undergo the pre-process, and feature extract. The data collected are cleaned by using Phyton to remove any outliers to ease the pre-processing process. These two processes are required in the development of machine learning models to obtain a clean dataset, which will make algorithms more convenient and accurate.

```python
import pandas as pd
from sklearn.impute import SimpleImputer
from sklearn.preprocessing import StandardScaler
from sklearn.feature_selection import SelectKBest
from sklearn.feature_selection import f_regression
```

Figure 5: Import necessary libraries

```
data = pd.read_csv('microclimate_data.csv')
```

Figure 6: Read the microclimate data from a CSV file

```
imputer = SimpleImputer(strategy='mean')
data[['temperature', 'relative_humidity', 'wind', 'solar_radiation', 'rainfall']] =
imputer.fit_transform
(data[['temperature', 'relative_humidity', 'wind', 'solar_radiation', 'rainfall']])

scaler = StandardScaler()
data[['temperature', 'relative_humidity', 'wind', 'solar_radiation', 'rainfall']] =
scaler.fit_transform
(data[['temperature', 'relative_humidity', 'wind', 'solar_radiation', 'rainfall']])
```

Figure 7: (Pre-processing) Handle missing values using mean imputation and
standardize the data

```
X = data[['temperature', 'relative_humidity', 'wind', 'solar_radiation', 'rainfall']]
y = data['target_variable']  # Target variable

kbest_selector = SelectKBest(score_func=f_regression, k=5)
X_new = kbest_selector.fit_transform(X, y)
selected_features = X.columns[kbest_selector.get_support()]
```

Figure 8: (Feature Engineering) Extract relevant features using SelectKBest with
f_regression as the scoring function

### 4.4.3 Splitting of Data into Training and Testing Sets

In the below code snippet, the train_test_split function from scikit-learn is used to split the pre-processed data (X_new) and the corresponding target variable (y) into training and testing sets. The test_size parameter is set to 0.2, indicating that 20% of the data will be used for testing, while the remaining 80% will be used for training. The random_state parameter is set to 42 to ensure reproducibility of the split. Then, it will proceed with training and evaluating machine learning models using the X_train, y_train for training, and X_test, y_test for testing.

```
from sklearn.model_selection import train_test_split

# Split the preprocessed data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X_new, y, test_size=0.2, random_state=42)
```

Figure 9: Split the pre-processed data into training and testing sets
with an 80-20 ratio

### 4.4.4 Training and Evaluation of Machine Learning Models

The next step after splitting the data into training and testing sets is to train and evaluate the machine learning models. The figure shows the code for training and evaluating the Random Forest and XGBoost algorithms: In this code snippet, the Random Forest and XGBoost models are initialized and trained using the training data (X_train and y_train). Then, predictions are made on the testing data (X_test) using the trained models. The mean absolute error (MAE) and mean squared error (MSE) are calculated to evaluate the performance of both models.

```
from sklearn.ensemble import RandomForestRegressor
from xgboost import XGBRegressor
from sklearn.metrics import mean_absolute_error, mean_squared_error

# Initialize and train the Random Forest model
rf_model = RandomForestRegressor()
rf_model.fit(X_train, y_train)

# Make predictions on the testing set using the trained Random Forest model
rf_predictions = rf_model.predict(X_test)

# Calculate evaluation metrics for the Random Forest model
rf_mae = mean_absolute_error(y_test, rf_predictions)
rf_mse = mean_squared_error(y_test, rf_predictions)

# Initialize and train the XGBoost model
xgb_model = XGBRegressor()
xgb_model.fit(X_train, y_train)

# Make predictions on the testing set using the trained XGBoost model
xgb_predictions = xgb_model.predict(X_test)

# Calculate evaluation metrics for the XGBoost model
xgb_mae = mean_absolute_error(y_test, xgb_predictions)
xgb_mse = mean_squared_error(y_test, xgb_predictions)

# Print the evaluation results
print("Random Forest - Mean Absolute Error:", rf_mae)
print("Random Forest - Mean Squared Error:", rf_mse)
print("XGBoost - Mean Absolute Error:", xgb_mae)
print("XGBoost - Mean Squared Error:", xgb_mse)
```

Figure 10: Train and evaluate the Random Forest and XGBoost algorithms

## 4.5    Parameter and Testing Methods

In order to evaluate the effectiveness of the proposed solution, several parameters are measured during the testing phase. These parameters provide valuable information about the performance and efficacy of the developed models. The parameters to be measured include:

### 4.5.1   Parameters to be Measured

#### 4.5.1.1 Prediction Accuracy

The accuracy of the predictive models in capturing and predicting the microclimate conditions is measured. This indicates how well the models can estimate the actual values of temperature, humidity, wind speed, solar radiation, and rainfall.

**4.5.1.2 Mean Absolute Error**

MAE is measured to determine the average absolute difference between the predicted and actual values. It provides insights into the average prediction error across all the microclimate parameters.

**4.5.1.3 Mean Squared Error**

MSE is calculated to assess the average squared difference between the predicted and actual values. It measures the overall variance between the predicted and actual values.

**4.5.2    Testing Procedure**

**4.5.2.1 Splitting Data**

The pre-processed microclimate data is divided into training and testing sets, typically using an 80:20 split. The training set is used to train the machine learning models, while the testing set is used for evaluating the performance.

**4.5.2.2 Model Evaluation**

The trained Random Forest and XGBoost models are applied to the testing set to make predictions for the microclimate parameters. The predicted values are then compared with the actual values from the testing set.

**4.5.2.3 Calculation of Evaluation Metrics**

The evaluation metrics, such as MAE and MSE, are calculated based on the predicted and actual values. These metrics provide quantitative measures of the model's performance.

**4.5.2.4 Analysis and Interpretation**

The evaluation results are analyzed to gain insights into the accuracy and effectiveness of the trained models. The performance of the models is assessed based on the evaluation metrics and compared to determine the superior algorithm for microclimate prediction.

**4.6    Chapter Summary**

In this chapter, the research experimental design and implementation of the proposed solution for preserving cultural heritage sites through microclimate monitoring and prediction are summarized. The steps for data collection, pre-processing, model development using Random Forest and XGBoost, evaluation and testing are outlined. The experimental setup and parameter details are provided, along with the evaluation metrics used to assess the performance of the models. The next chapter will present the results and analysis of the experiments, highlighting the insights gained and the recommendations for preserving cultural heritage sites.

# CHAPTER 5

## CONCLUSION AND RECOMMENDATIONS

### 5.1    Research Outcomes

This chapter presents the outcomes of the research conducted on preserving cultural heritage sites through the application of the Random Forest and XGBoost algorithms for microclimate monitoring and prediction. The research aimed to develop an effective and efficient approach to monitor and predict the microclimate conditions at cultural heritage sites, thereby aiding in the preservation of these sites for future generations.

One of the key research outcomes is the development of a microclimate monitoring system utilizing the Random Forest and XGBoost algorithms. The system collects real-time data from various sensors deployed at cultural heritage sites, including temperature, humidity, light intensity, and air quality. The collected data is then processed and analyzed using the Random Forest and XGBoost algorithms to identify patterns and trends in microclimate conditions.

Another significant outcome of this research is the achievement of accurate microclimate prediction at cultural heritage sites. By training the Random Forest and XGBoost algorithms on historical microclimate data, the developed system can forecast future microclimate conditions with a high degree of accuracy. This prediction capability enables heritage site managers and conservationists to proactively plan and implement appropriate preservation strategies based on anticipated changes in the microclimate.

**5.2     Contributions to Knowledge**

The research conducted in this thesis has made several contributions to the field of cultural heritage preservation and microclimate monitoring. These contributions include the application of machine learning algorithms, specifically the Random Forest and XGBoost algorithms, in the context of microclimate monitoring and prediction at cultural heritage sites. By demonstrating the effectiveness of these algorithms in capturing complex relationships between various environmental factors and microclimate conditions, this research provides valuable insights into the potential of machine learning techniques for heritage site preservation.

The research presents a comprehensive framework for monitoring and predicting microclimate conditions at cultural heritage sites. This framework integrates data collection, preprocessing, analysis, and prediction using the Random Forest and XGBoost algorithms. The developed framework can serve as a guide for future researchers and practitioners in the field of cultural heritage preservation, providing a structured approach to leveraging machine learning for microclimate management.

By accurately monitoring and predicting microclimate conditions, this research contributes to the development of improved preservation strategies for cultural heritage sites. The insights gained from the analysis of microclimate data can inform decision-making processes related to site maintenance, climate control, and artifact preservation, ultimately enhancing the long-term sustainability of these important cultural assets.

## 5.3    Future Works

While this research has achieved significant milestones in the preservation of cultural heritage sites through microclimate monitoring and prediction, there are several avenues for future research and development. Some potential areas of focus include the integration of additional data sources, such as weather forecasts, aerial imagery, and historical records, to further enhance the accuracy and reliability of microclimate prediction models. Incorporating these diverse data sets can provide a more comprehensive understanding of the factors influencing microclimate conditions and enable more robust decision-making processes.

Continued monitoring and analysis of microclimate conditions over extended periods can yield valuable insights into the long-term trends and impacts on cultural heritage sites. Future research should consider conducting longitudinal studies to capture the dynamic nature of microclimate conditions and evaluate the effectiveness of preservation strategies over time.

Encouraging collaboration and knowledge sharing among researchers, practitioners, and stakeholders in the field of cultural heritage preservation is crucial for advancing the application of microclimate monitoring and prediction techniques. Future research should focus on establishing platforms for collaboration, fostering interdisciplinary partnerships, and promoting the dissemination of research findings to maximize the impact on heritage site conservation efforts.

By addressing these future research directions, the field of microclimate monitoring and prediction for cultural heritage preservation can continue to evolve and contribute to the sustainable management of these invaluable cultural assets.

**REFERENCES**

Lombardo, L., Tanyas, H., &amp; Nicu, I. C. (2020). Spatial modeling of multi-hazard threat to Cultural Heritage Sites. Engineering Geology, 277, 105776. https://doi.org/10.1016/j.enggeo.2020.105776

Fabbri, K., &amp; Bonora, A. (2021). Two new indices for preventive conservation of the cultural heritage: Predicted risk of damage and Heritage Microclimate Risk. Journal of Cultural Heritage, 47, 208–217. https://doi.org/10.1016/j.culher.2020.09.006

Pioppi, B., Pigliautile, I., Piselli, C., &amp; Pisello, A. L. (2020). Cultural Heritage Microclimate Change: Human-centric approach to experimentally investigate intra-urban overheating and numerically assess foreseen future scenarios impact. Science of The Total Environment, 703, 134448. https://doi.org/10.1016/j.scitotenv.2019.134448

Alcaraz Tarragüel, A., Krol, B., &amp; van Westen, C. (2012). Analysing the possible impact of landslides and avalanches on cultural heritage in Upper Svaneti, Georgia. Journal of Cultural Heritage, 13(4), 453–461. https://doi.org/10.1016/j.culher.2012.01.012

Sevetlidis, V., &amp; Pavlidis, G. (2019). Effective raman spectra identification with tree-based methods. Journal of Cultural Heritage, 37, 121–128. https://doi.org/10.1016/j.culher.2018.10.016

Kobayashi, K., Hwang, S.-W., Okochi, T., Lee, W.-H., &amp; Sugiyama, J. (2019). Non-destructive method for wood identification using conventional X-ray computed tomography data. Journal of Cultural Heritage, 38, 88–93. https://doi.org/10.1016/j.culher.2019.02.001

Zou, Z., Zhao, X., Zhao, P., Qi, F., &amp; Wang, N. (2019). CNN-based statistics and location estimation of missing components in routine inspection of Historic Buildings. Journal of Cultural Heritage, 38, 221–230. https://doi.org/10.1016/j.culher.2019.02.002

CC Publications Online. ICOM. (n.d.). https://www.icom-cc-publications-online.org/4417/Teaching-machines-to-think-like-conservators--Machine-

learning-as-a-tool-for-predicting-the-stability-of-paper-based-archive-and-library-collections

Kejser, U. B., Ryhl-Svendsen, M., Boesgaard, C., &amp; Hansen, B. V. (n.d.). Teaching machines to think like conservators – Machine learning as a tool for predicting the stability of paper-based archive and library collections. Transcending Boundaries: Integrated Approaches to Conservation. ICOM-CC 19th Triennial Conference Preprints, Beijing, 17–21 May 2021. https://www.icom-cc-publications-online.org/4417/Teaching-machines-to-think-like-conservators--Machine-learning-as-a-tool-for-predicting-the-stability-of-paper-based-archive-and-library-collections

Pei, J., Gong, J., &amp; Wang, Z. (2020). Risk prediction of household mite infestation based on machine learning. Building and Environment, 183, 107154. https://doi.org/10.1016/j.buildenv.2020.107154

Lowenthal, D. (2015). The Past is a Foreign Country. Cambridge: Cambridge University Press.

UNESCO. (1972). Convention Concerning the Protection of the World Cultural and Natural Heritage. Paris: UNESCO.

Smith, L. (2006). Uses of Heritage. London: Routledge.

Timothy, D. J., & Boyd, S. W. (2003). Heritage Tourism. Harlow: Prentice Hall.

Cassar, M. (2005). Climate Change and the Historic Environment. London: English Heritage.

Camuffo, D. (2014). Microclimate for Cultural Heritage: Conservation, Restoration, and Maintenance of Indoor and Outdoor Monuments. Amsterdam: Elsevier.

Lankester, P., & Brimblecombe, P. (2012). The impact of future climate change on historic interiors. Science of The Total Environment, 417-418, 248-254.

Yu, T., Lin, C., Zhang, S., Wang, C., Ding, X., An, H., Liu, X., Qu, T., Wan, L., You, S., Wu, J., &amp; Zhang, J. (2022). Artificial Intelligence for Dunhuang Cultural Heritage Protection: The project and the dataset. International Journal of Computer Vision, 130(11), 2646–2673. https://doi.org/10.1007/s11263-022-01665-x

Kumar, P., Ofli, F., Imran, M., &amp; Castillo, C. (2020). Detection of disaster-affected cultural heritage sites from social media images using Deep Learning Techniques. Journal on Computing and Cultural Heritage, 13(3), 1–31. https://doi.org/10.1145/3383314

Staniforth, S. (2013). Historical Perspectives on Preventive Conservation. Getty Conservation Institute.

Stovel, H., Stanley-Price, N., &amp; Killick, R. G. (2005). Conservation of living religious heritage: Papers from the ICCROM 2003 forum on living religious history: Conserving the sacred. International Centre for the Study of the Preservation and Restoration of Cultural Property.

Muñoz Viñas (2002) Contemporary theory of conservation, Studies in Conservation, 47:sup1, 25-34, DOI: 10.1179/sic.2002.47.Supplement-1.25

Ashley-Smith, J. (2016).Risk assessment for object conservation.  Routledge.

Caple, C. (2012). Preventive conservation in museums. Routledge.

Matero, F. (1999). Lessons from the Great House: Condition and treatment history as prologue to site conservation and management at Casa Grande Ruins National Monument. Conservation and Management of Archaeological Sites, 3(4), 203–224. https://doi.org/10.1179/135050399793138482

Prieto, A. J., Silva, A., de Brito, J., Macías-Bernal, J. M., &amp; Alejandre, F. J. (2017). Multiple linear regression and fuzzy logic models applied to the functional service life prediction of Cultural Heritage. Journal of Cultural Heritage, 27, 20–35. https://doi.org/10.1016/j.culher.2017.03.004

Gonthier, N., Gousseau, Y., Ladjal, S., &amp; Bonfait, O. (2019). Weakly supervised object detection in artworks. Lecture Notes in Computer Science, 692–709. https://doi.org/10.1007/978-3-030-11012-3_53

Valero, E., Forster, A., Bosché, F., Hyslop, E., Wilson, L., &amp; Turmel, A. (2019). Automated defect detection and classification in ashlar masonry walls using machine learning. Automation in Construction, 106, 102846. https://doi.org/10.1016/j.autcon.2019.102846

APPENDIX A Gantt Chart for FYP 1

| PHASE/WEEK | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **PLANNING PHASE** | | | | | | | | | | | | | | | | | | | | | | | |
| Research Supervisor and Topic Research Selection | █ | █ | | | | | | | | | | | | | | | | | | | | | |
| Research Proposal Documentation and Submission | | | █ | | | | | | | | | | | | | | | | | | | | |
| Research Proposal Interview | | | █ | | | | | | | | | | | | | | | | | | | | |
| Research Proposal Correction | | | █ | █ | | | | | | | | | | | | | | | | | | | |
| **ANALYSIS PHASE** | | | | | | | | | | | | | | | | | | | | | | | |
| Chapter 1 Introduction Documentation | | | | █ | █ | | | | | | | | | | | | | | | | | | |
| Collection of research source, material and information searching | | | | █ | █ | █ | | | | | | | | | | | | | | | | | |
| Chapter 2 Literature Review Documentation | | | | | | █ | █ | | | | | | | | | | | | | | | | |
| Analysis of existing research | | | | | | | | █ | █ | | | | | | | | | | | | | | |
| Chapter 3 Research Methodology Documentation | | | | | | | | | █ | █ | █ | | | | | | | | | | | | |
| **RESULT AND FINDINGS PHASE** | | | | | | | | | | | | | | | | | | | | | | | |
| Explore solution approach | | | | | | | | | | | █ | █ | | | | | | | | | | | |
| Information gathering on existing application analysis | | | | | | | | | | | | █ | █ | | | | | | | | | | |
| Chapter 4 Research and Design Implementation Documentation | | | | | | | | | | | | | | █ | █ | █ | | | | | | | |
| Chapter 5 Conclusion Documentation | | | | | | | | | | | | | | | | | █ | | | | | | |
| Report Compilation | | | | | | | | | | | | | | | | | | █ | █ | █ | | | |
| **FYP 1** | | | | | | | | | | | | | | | | | | | | | | | |
| Presentation of FYP1 | | | | | | | | | | | | | | | | | | | | | █ | | |
| Report Correction | | | | | | | | | | | | | | | | | | | | | █ | █ | █ |

44

APPENDIX B Gantt Chart for FYP 2

| PHASE/WEEK | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **FYP 2** | | | | | | | | | | | | | | | | | | | |
| FYP 2 Briefing | ■ | | | | | | | | | | | | | | | | | | |
| **ANALYSIS PHASE** | | | | | | | | | | | | | | | | | | | |
| Chapter 1 Introduction Review | ■ | | | | | | | | | | | | | | | | | | |
| Collection of research source, material and information searching | | ■ | | | | | | | | | | | | | | | | | |
| Chapter 2 Literature Review | | | ■ | | | | | | | | | | | | | | | | |
| Analysis of existing research | | | | ■ | | | | | | | | | | | | | | | |
| Chapter 3 Research Methodology Review | | | | | ■ | ■ | | | | | | | | | | | | | |
| **RESULT AND FINDINGS PHASE** | | | | | | | | | | | | | | | | | | | |
| Explore solution approach | | | | | | | ■ | | | | | | | | | | | | |
| Information gathering on existing application analysis | | | | | | | | ■ | | | | | | | | | | | |
| Chapter 4 Machine Learning Model Development | | | | | | | | | ■ | | | | | | | | | | |
| Experiment | | | | | | | | | | ■ | ■ | | | | | | | | |
| Chapter 5 Dashboard Development | | | | | | | | | | | | ■ | ■ | | | | | | |
| Chapter 6 Conclusion Documentation and Review | | | | | | | | | | | | | ■ | ■ | | | | | |
| Thesis Compilation | | | | | | | | | | | | | | | | | | | |
| **FYP 2** | | | | | | | | | | | | | | | ■ | | | | |
| Presentation of FYP1 | | | | | | | | | | | | | | | ■ | ■ | | | |
| Thesis Correction | | | | | | | | | | | | | | | | | | ■ | ■ |