

FINAL YEAR PROJECT

Text Classification on Diabetes Mellitus
Symptom and Treatment Documents Using
Support Vector Machine (SVM)

Presentation Video:

https://youtu.be/RvpUL_deGNU

Presented by

Phang Cheng Yi

Supervised by

Dr Sharin Hazlin Binti Huspi

Dr Ahmad Najmi Bin Amerhaider Nuar



Table of Contents



1

CHAPTER 1 : INTRODUCTION



2

CHAPTER 2 : LITERATURE REVIEW



3

CHAPTER 3 : RESEARCH METHODOLOGY



4

CHAPTER 4 : RESEARCH DESIGN AND IMPLEMENTATION



5

CHAPTER 5 : CONCLUSION AND RECOMMENDATIONS





CHAPTER 1

INTRODUCTION

Introduction

■ Diabetes Mellitus (DM)

- Metabolic disorder characterised by excessively increased blood glucose levels with numerous subtypes
- Severe and varied symptoms of hyperglycemia include abnormalities in the metabolism of carbohydrates, fats, and proteins
- Early detection and efficient management is required but challenging

■ Text Classification

- A type of the NLP unstructured text analysis techniques
- A predetermined label or tag will be given to each document in the dataset by the classifier

■ Support Vector Machine (SVM)

- A type of commonly used machine learning text classifiers that give significant result

Problem Background



The number of people suffered from the disease Diabetes Mellitus (DM) keep increasing

Early diagnosis by identifying the symptoms is significant to ascertain suitable treatments

Difficulty to discover and to classify the important information from numerous documents for better understanding and is time consuming

A lot of publications regarding to diabetes for early diagnosis, treatment, management and prevention

Previous studies mostly used clinical data and patient medical record in classification using machine learning methods such as Fine Decision Tree and SVM but lack of study focus on the classification for medical journal articles that describing the symptoms and treatments of DM

Potential to develop text classification model for diabetes symptom and treatment documents using SVM method to optimize the diabetes diagnosis and management

Problem Statement

- Overwhelming amount of medical literature and research on DM, which can hinder the efficiency of early detection and the effectiveness of management for diabetes patients and doctors
- Keeping up with the latest research discoveries and therapeutic approaches is getting harder as diabetes becomes more common and more complex
- Lack of research on the application of SVM for text classification on DM symptom and treatment documents
- Text classification model is used to assist in finding crucial symptoms and methods of therapy for DM, enhancing the effectiveness and precision of information retrieval



Research Objectives

Goal

To develop and evaluate a SVM based text classification model for DM symptom and treatment documents

Objectives

(a) To identify the related features that are relevant to Diabetes Mellitus (DM) symptoms and treatments in multiple documents.

(b) To perform text classification for a collection of DM documents dataset using Support Vector Machine (SVM) method.

(c) To evaluate the performance of machine learning model that apply Support Vector Machine (SVM) method through several model evaluation techniques

Research Scope

1 To use a collection of articles from PubMed by National Center for Biotechnology Information (NCBI)

2 The documents in the dataset are chosen by focusing on the symptoms and treatments of DM

3 To use Term Frequency-Inverse Document Frequency (TF-IDF) algorithms to spot the important terms of DM symptoms and treatments

4 To use Support Vector Machine (SVM) algorithms to classify text documents



Research Contribution

■ First

Contribute to the fields of natural language processing (NLP) and machine learning

■ Second

Contribute to the information retrieval fields with the enhancement of retrieving process that enable domain stakeholders of DM to efficiently identify and classify the important symptoms and treatments for DM based on the analysis of a large corpus of medical documents

■ Third

Potentially facilitate the development of more effective interventions for diabetes management.



CHAPTER 2

LITERATURE REVIEW

Diabetes Mellitus (DM)

Four major types of DM:

- Type 1 diabetes (T1DM)
- Type 2 diabetes (T2DM)
- Gestational diabetes mellitus (GDM)
- Other specific types

Misdiagnosis always occur between T1DM and T2DM. Although they share some certain features, but the way of presenting the symptoms is distinct.

- T1DM: Quickly within few weeks
- T2DM: Slowly over a long period of time

Can cause many vital organ fail to function, e.g. retina, kidney, neurological system, heart, blood vessels

Formed when the body's cells and tissues are **unable to use the insulin produced by the pancreas** or when the **pancreas cannot produce enough insulin**

A disease that involve metabolic disorder distinguished by **hyperglycemia**, a physiologically dysfunctional condition reflected by **excessively increased blood glucose levels**

Diabetic with T1DM and T2DM always accompanied by the following **common symptoms**:

weight loss, dry or itchy skin, blurred vision, polyuria, excessive thirst, increased hunger, sluggish wound healing, frequent infections and depressive mood

Two main drug therapies in **treating people with T2DM**:

- Oral
- Injection

Therapies in **treating people with T1DM**:

- Insulin with adjunctive drugs like Metformin and DPP-4 Inhibitors

Certain cases of T1DM and T2DM:

- Combination of multiple drugs
- Insulin for T2DM diabetics

Common treatment:

- Diet monitoring



Text Mining

Interdisciplinary field: Information retrieval, text analysis, information extraction, categorization, clustering, visualisation, data mining, and machine learning.

A data mining approach or process that identifies previously unexplored and insightful information from massive quantities of unstructured text data

Overall process of text mining:

- Problem Identification
- Text pre-processing
- Pattern extraction (Text mining)
- Text post-processing (Evaluation)
- Knowledge usage



Text Classification

- The automated assignment of text documents into predefined groups according to the content of the text itself by using various technologies and algorithms
- One the most prominent technique of text mining and Natural Language Processing (NLP)
- Can be break down into four stages



Stages of Text Classification System



Stage 1

Feature Extraction

Stage 2

Dimension Reduction

Stage 3

Classifier Selection

Stage 4

Evaluation

Text Enrichment

Word2Vec

- One of the word embedding methods
- Uses deep learning ideology to create a distributed representation of words in semantic space where a corresponding vector will be assigned to every term
- To create a distributed representation of the word:
 - Continuous skip-gram or Continuous Bag-of-Words (CBOW)
- To determine words' similarity:
 - Euclidean distance, Cosine Similarity
- **Benefits:**
 - Excellent training effectiveness and rich semantics, which can be utilised for part-of-speech analysis, **synonym search**, and clustering

Feature Extraction

- Convert text into keyword schedule to makes it feasible for supervised learning
- Vectorizing text data is necessary before utilising it as input for a machine learning system since the input must be numerical
- Two main types of feature extraction methods:
 - Weighted word
 - Bag-of-Words (BoW), TF-IDF
 - Word embeddings
 - Word2Vec, GloVe



Weighted Word Technique

TF-IDF

- Emphasises a word's importance in a text in relation to the entire corpus
- Basic concept:
 - a term that frequently appears in a document but infrequently in the whole corpus is more informative than a word that regularly appears in both corpus and document
- Two steps involved in TF-IDF, calculating:
 - Term Frequency (TF)
 - Inverse Document Frequency (IDF)

Summary of Machine Learning Approaches

Approaches	Advantages	Limitations
Support Vector Machine (SVM)	<ul style="list-style-type: none">• Efficient in performing non-linear classification.• Low prediction error due to risk minimization concept.• Great option for handling high-dimensional data.	<ul style="list-style-type: none">• Limited to small sample.• Selection of different kernels may yield the efficiency of classifier.
Naive Bayes Classifier (NBC)	<ul style="list-style-type: none">• Easy and quick to implement.	<ul style="list-style-type: none">• Precision might decrease when data size is small.• Correlated attributes will affect the performance.

Summary of Machine Learning Approaches

Approaches	Advantages	Limitations
K-Nearest Neighbor (KNN)	<ul style="list-style-type: none">• Simple to run.• Effective with big training datasets and robust to noisy data.	<ul style="list-style-type: none">• Slow in handling large data.• Sensitive to unnecessary factors.• High demand in memory.
Decision Tree	<ul style="list-style-type: none">• Data normalization is not necessary.• Support high dimension data.	<ul style="list-style-type: none">• The continuous fields are tough to forecast.• Unstable due to its dependency on dataset type.
Random Forest	<ul style="list-style-type: none">• High accuracy	<ul style="list-style-type: none">• Time consuming and large memory space needed.• Overfitting may occur if noise exist.

Summary of Related Work

Text Classification

- Several studies about text classification were conducted to compare different types of machine learning algorithms
- Datasets: Journal articles, Research paper, News articles, Movie Reviews and SMS messages

Evaluations

- SVM is used almostly in all the studies
- By comparing with other machine learning algorithms, SVM model always shows better performance
 - (Rasheed et al., 2018):
 - Accuracy: **SVM-68.73%**, Decision Tree-62.37%, KNN-55.41%
 - (Chowdhury and Schoen, 2020):
 - Accuracy: **SVM-88%**, Naive Bayes-86%, KNN-83%

Discussion

Clinical dataset (Pima Indian) and eye fundus numerical images dataset were commonly used among the researchers to classify whether the patients are diabetic or non-diabetic and to discover the cause of the diabetes.

Numerous publications pertinent to the most thoroughly studied diseases are readily available. These unstructured data have become the potential sources that can provide beneficial information to the medical experts in their diagnosis.

There is lack of studies that occupy thoughts with the narrative documents of DM by classifying their symptoms for different types of DM and its treatment.

SVM is able to perform well in text classification.

It is potential to conduct research studies on text classification for DM symptom and treatment documents using SVM algorithms

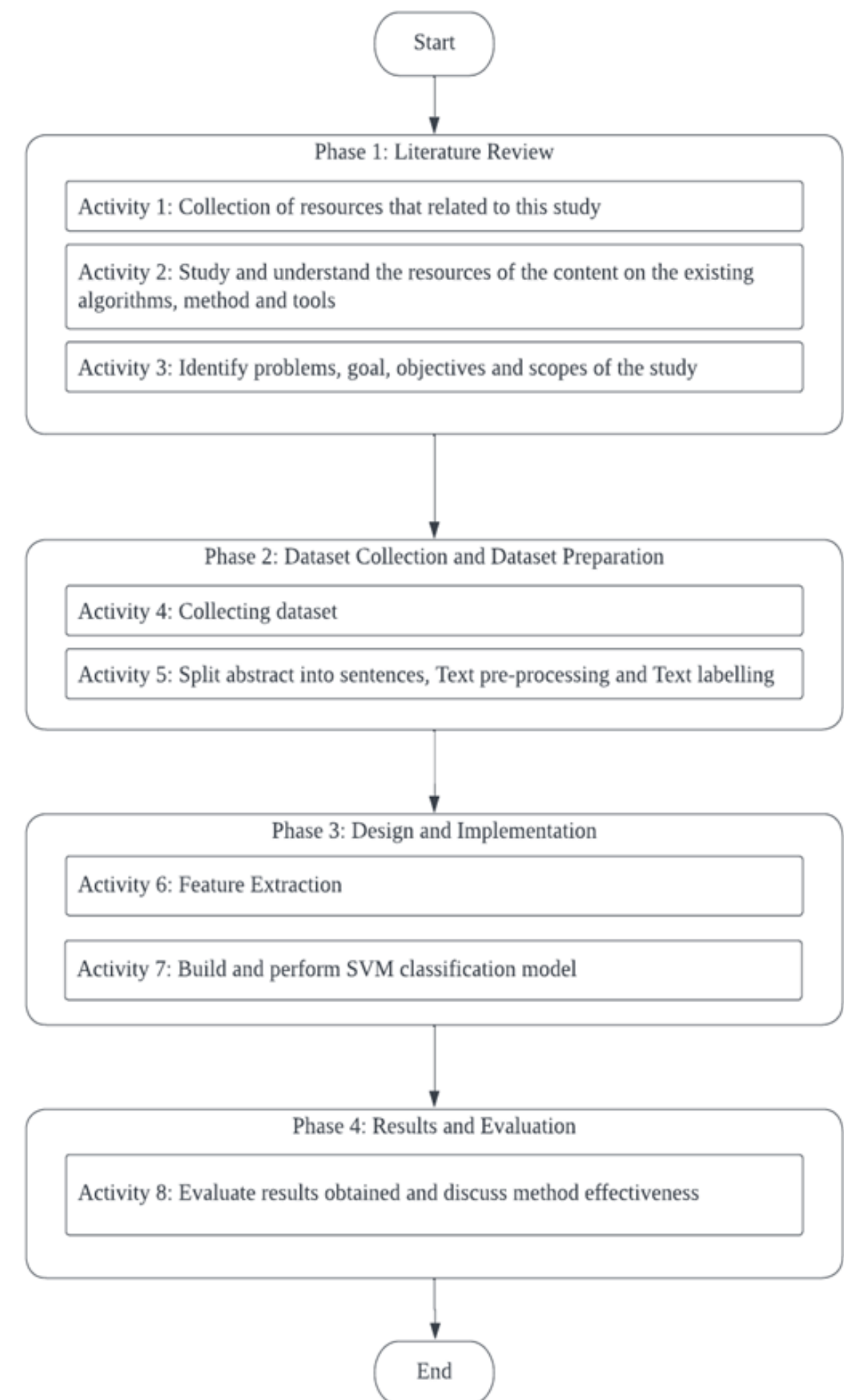


CHAPTER 3

RESEARCH

METHODOLOGIES

Research Workflow

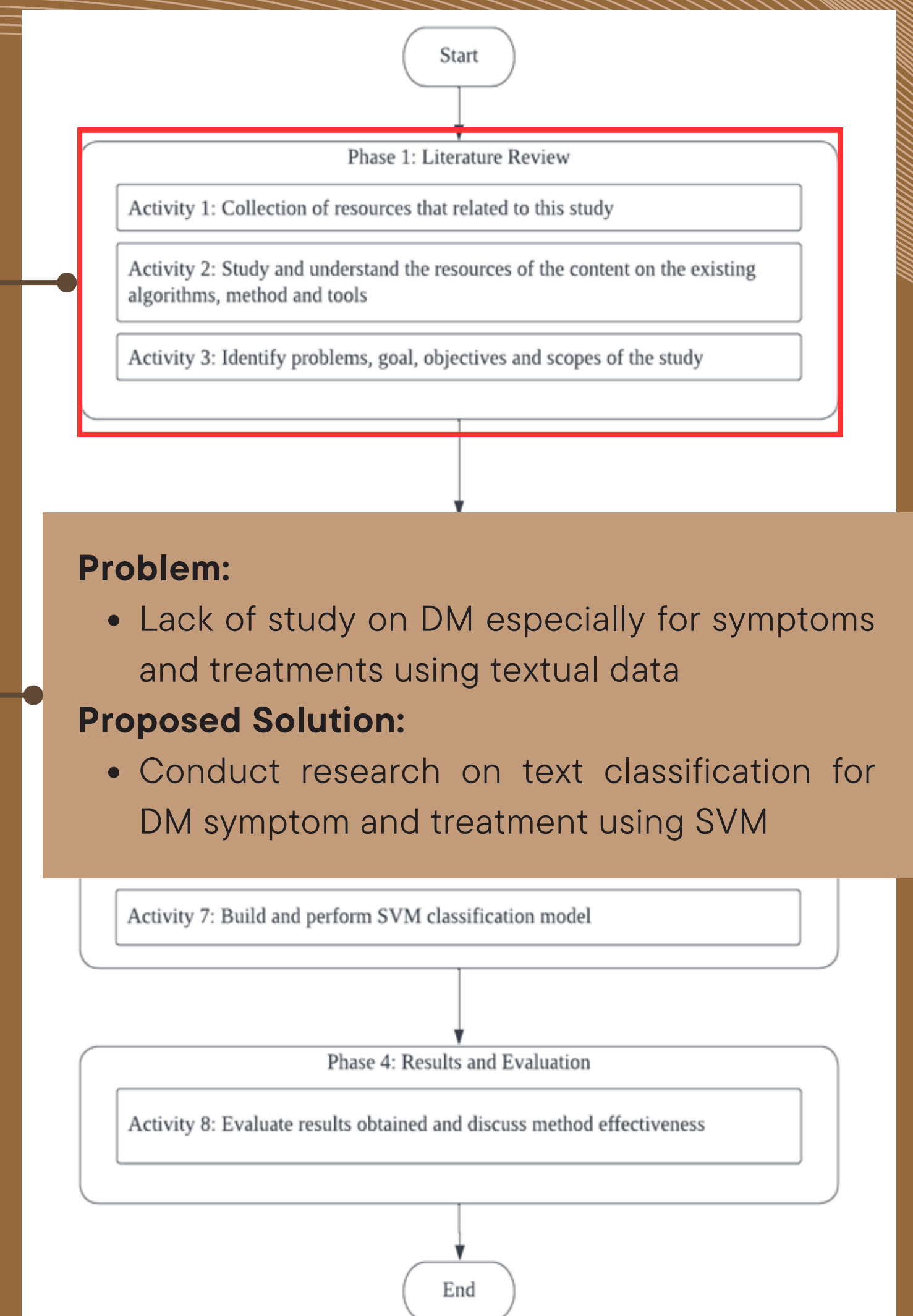


Research Workflow

Phase 1

- The problem background, past studies and related work are discussed.
- Have more understanding regarding to the definition and existing knowledge of DM, text classification and SVM.
- Previous studies: text classification using several types of machine learning methods
 - To find out how previous researchers did to conduct their research

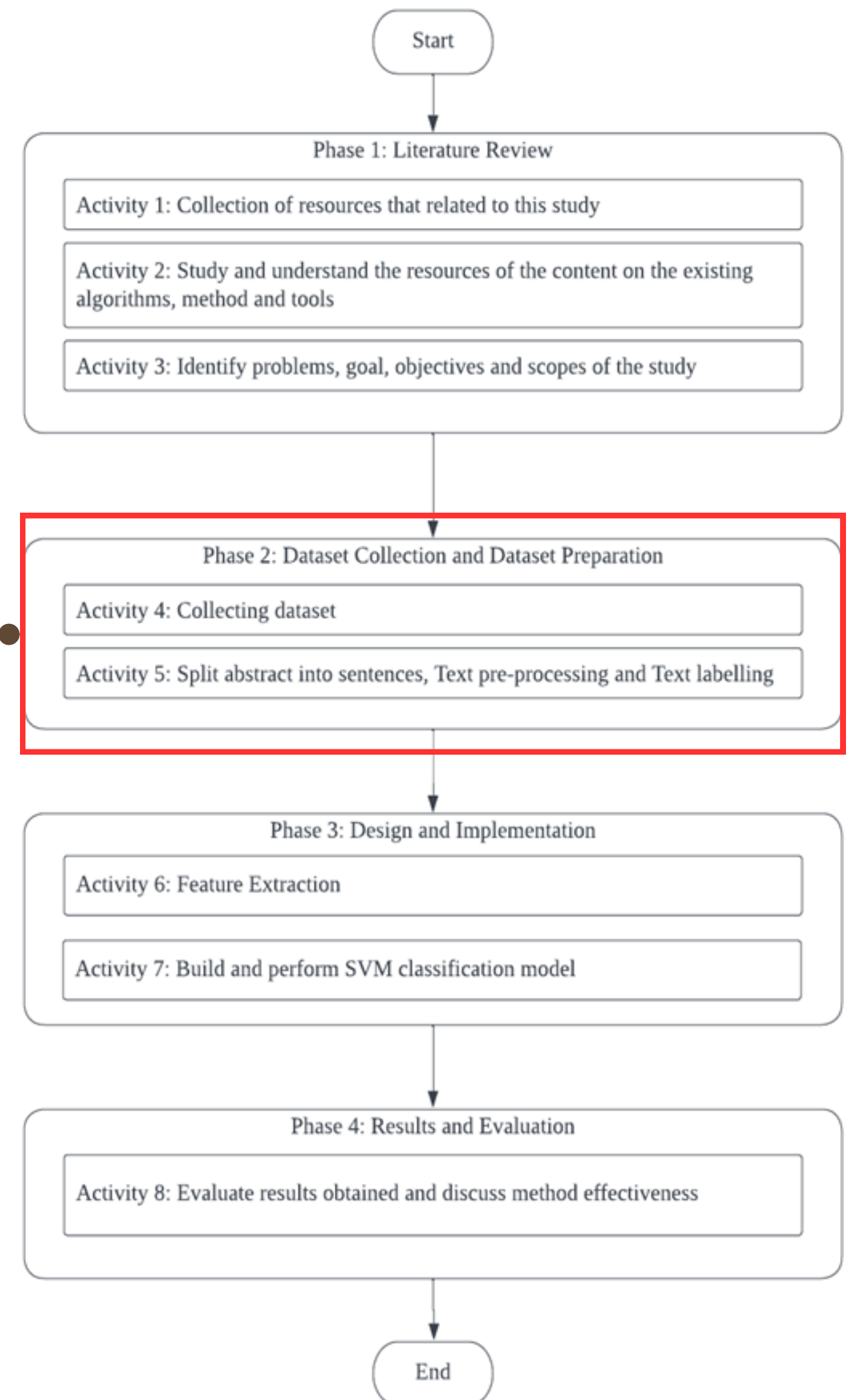
To comprehend the underlying issue or potential future work



Research Workflow

Phase 2

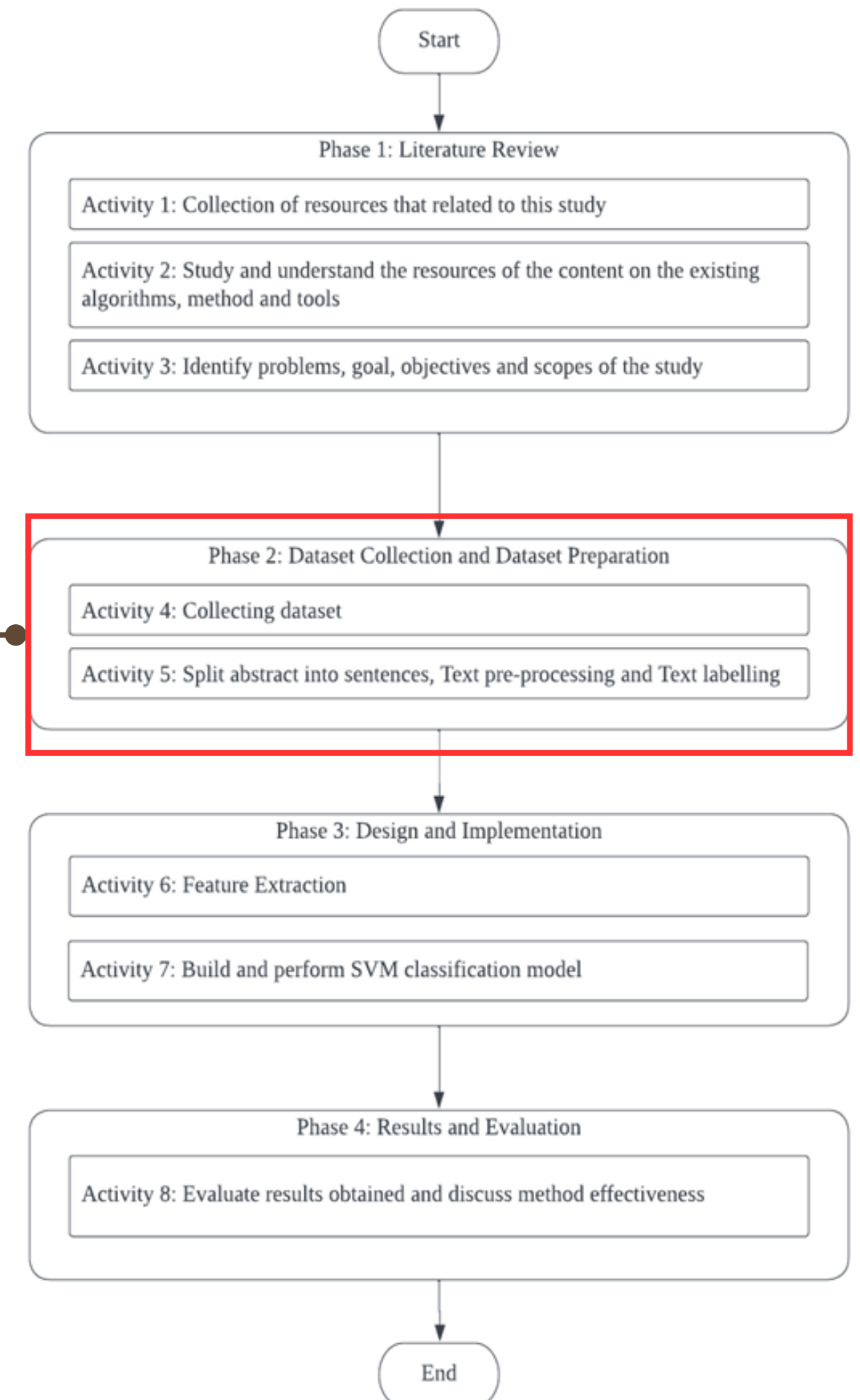
- Data collection:
 - **Source: PubMed website**
 - mainly accesses MEDLINE database of references and abstract on topics related to life science and biomedical.
 - allow user to look for articles using various criteria
 - The **search key term, range of year** and **size of dataset** is set to fetch the relevant journal articles.
 - Saved as csv file for better readability and easier processing



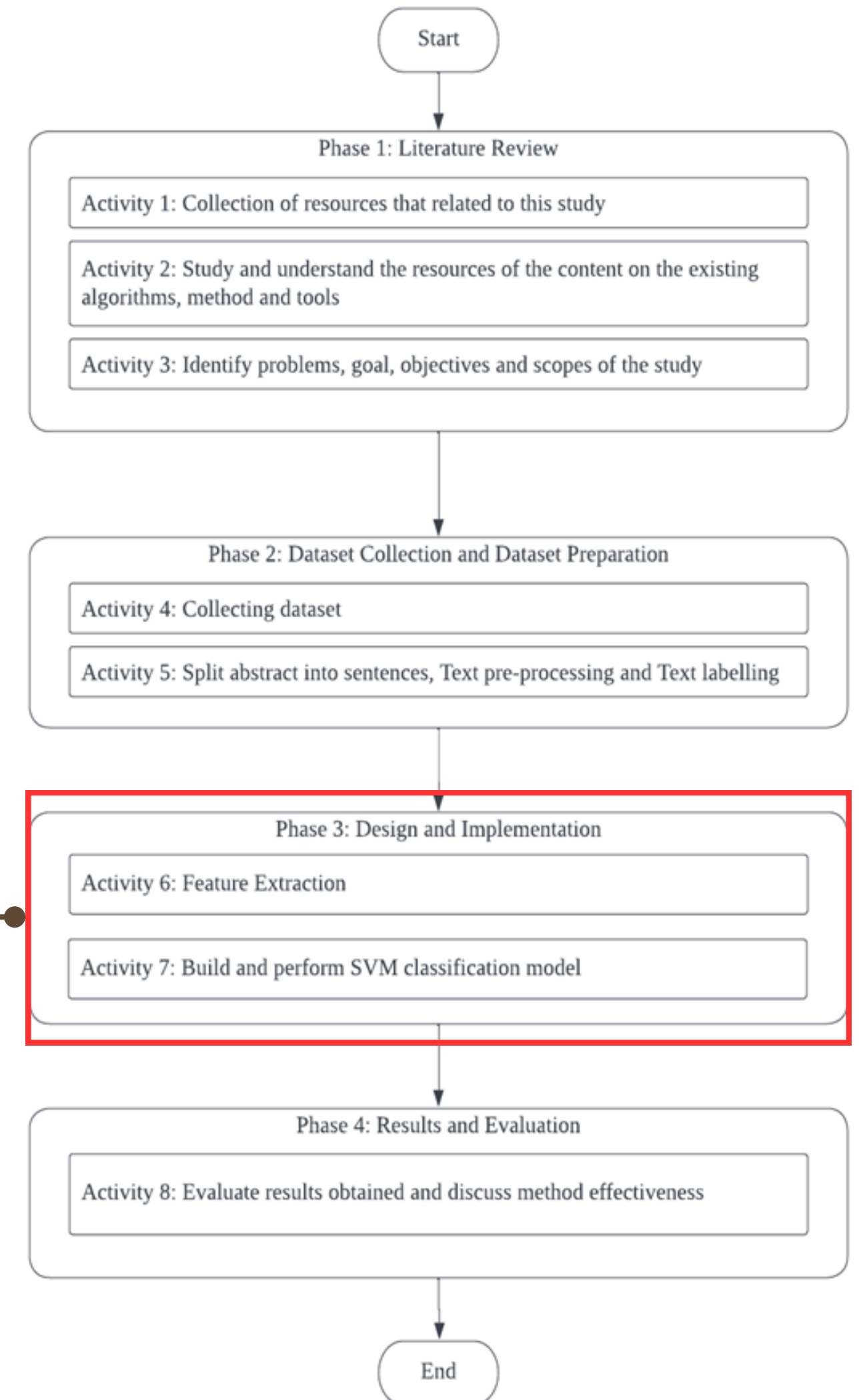
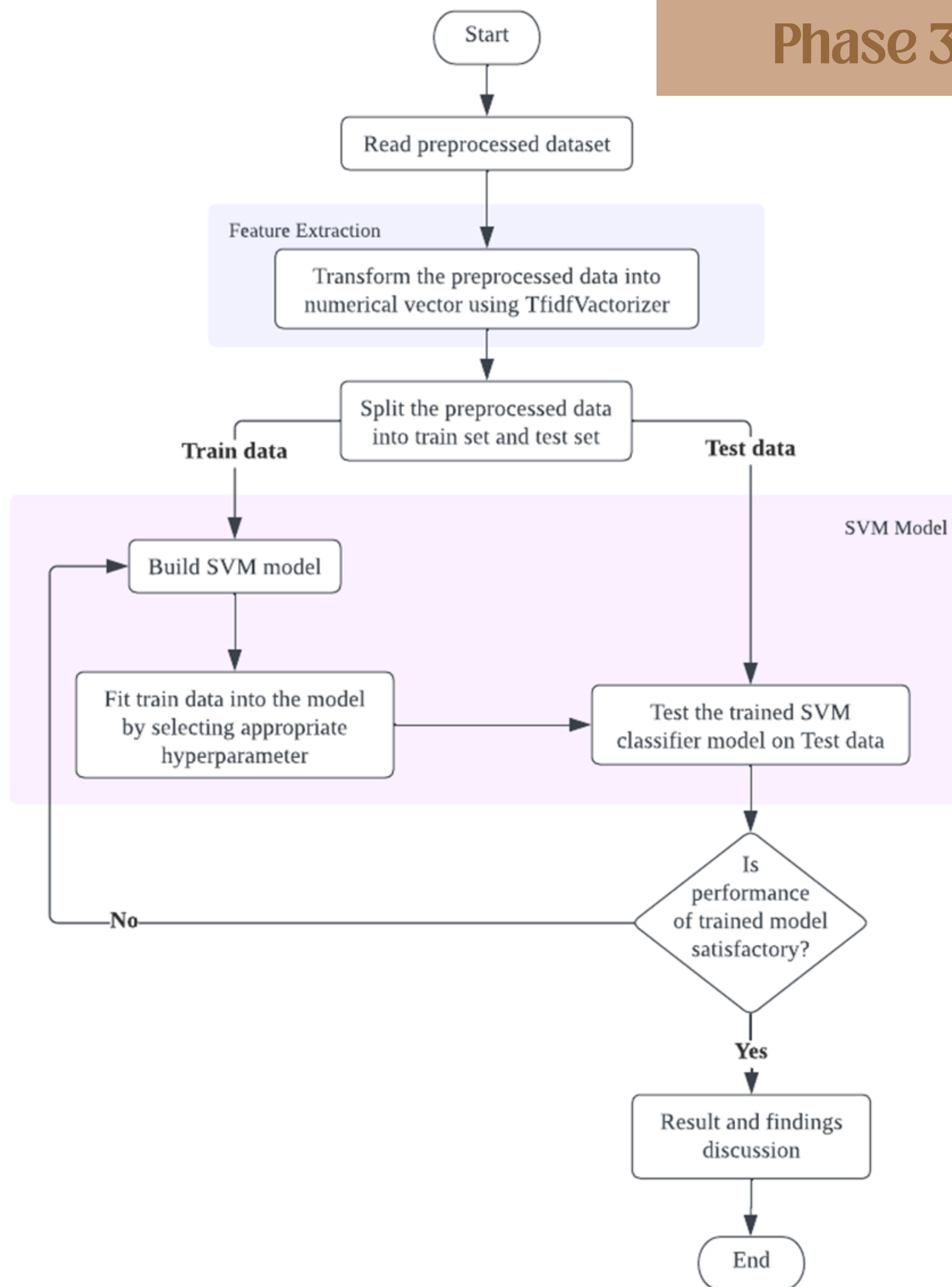
Research Workflow

Phase 2

- Dataset Preparation:
 - Split abstract into sentences
 - Text Pre-processing
 - Text Labelling
- Exploratory Data Analysis
 - Wordcloud



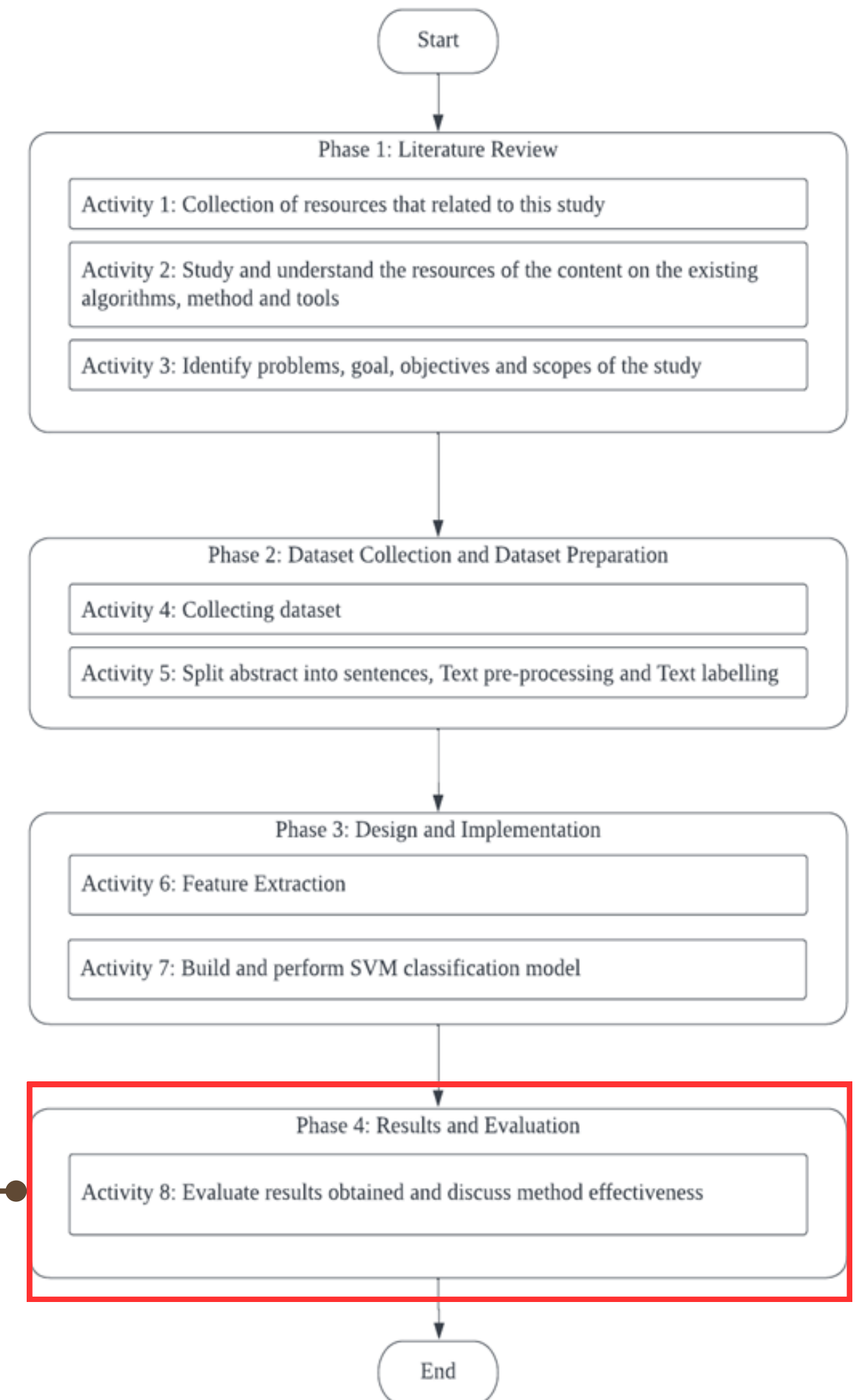
Phase 3



Research Workflow

Phase 4

- **Results evaluation and analysis** by comparing the performance when:
 - different kernel of SVM is used
 - different train test split ratio is used
- **Performance measurement:**
 - Confusion matrix
 - True Positive (TP)
 - True Negative (TN)
 - False Positive (FP)
 - False Negative (FN)
 - Metrics:
 - Accuracy
 - Precision
 - Recall





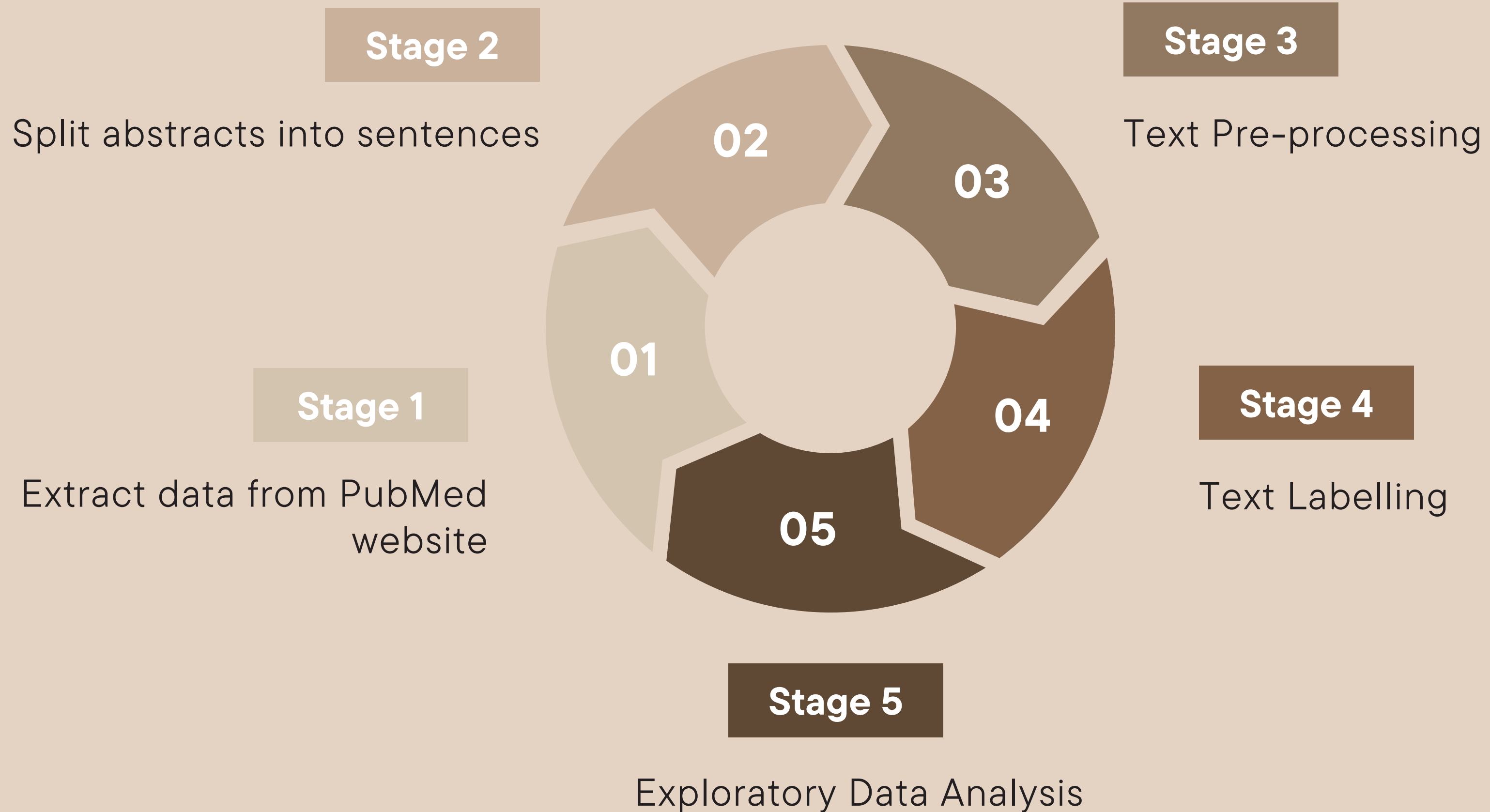
CHAPTER 4

RESEARCH DESIGN

AND

IMPLEMENTATION

Dataset Collection and Data Preparation



STAGE 1

Extract Data

- Dataset scrapped from PubMed website
- The search terms are defined as "Type 1 Diabetes Mellitus", "Type 2 Diabetes Mellitus", "Symptom" and "Treatment."
- The range of years is set from 2020 to 2023.
- The number of journal articles selected is 300.
- Collected data are saved as csv file.

STAGE 1

Pmid	Title	Abstract
37277527	The genetic architecture of type 2 diabetes	[Diabetes is one of the most common phenotypes of Wolfram syndrome owing to the presence of the variants of the WFS gene.]
37264885	[Not Available]	[Maturity-onset diabetes of the young (MODY) is a group of hereditary monogenetic forms of diabetes. MODY accounts for approximately 1-5% of all cases of diabetes mellitus.]
37264478	Higher β -cell function predicts better glycaemic control in type 2 diabetes	[Diabetes is a metabolic disorder of glucose homeostasis in which β -cell destruction occurs silently and is detected mainly by hyperglycaemia.]
37259043	Validation of the Diabetes Self-Rating Scale (DSRS)	[StringElement('Disordered eating behaviours (DEBs) in patients with type 1 diabetes mellitus (T1DM) are associated with altered eating patterns and weight changes.)]
37258468	SGLT2i in Type 2 Diabetes Mellitus	[Sodium-glucose cotransporter inhibitors (SGLT2i) play an increasingly important role in type 2 diabetes mellitus (T2DM) due to their ability to improve glycaemic control and reduce cardiovascular risk.]
37257909	Navigating the Seas of Glycemic Control: The Role of Continuous Glucose Monitoring in Type 1 Diabetes Mellitus.	
37241059	Benefit or Burden? Dapagliflozin in Type 2 Diabetes Mellitus	[Purpose: Dapagliflozin has been used extensively in patients with type 2 diabetes mellitus (T2DM). However, due to its mechanism of action, it may have both benefits and risks.]
37187225	Determination of the correlation between %TIR and HbA1c in pregnant women with type 1 diabetes mellitus	[To determine the correlation between %TIR and HbA1c in pregnant women with type 1 diabetes mellitus (D1DM).]
37161897	The Possibility of Treating Type 2 Diabetes Mellitus with Metformin	[This minireview discusses the very important biomedical problem of treating type 2 diabetes mellitus (T2D). T2D accounts for approximately 90% of all cases of diabetes mellitus.]
37133646	Use of Flash Glucose Monitoring to Assess Glycemic Control in Patients With Type 2 Diabetes Mellitus	[StringElement('Estimation of laboratory-derived glycated hemoglobin (HbA1c) cannot be individually used to monitor clinical status.)]
37131168	Understanding the Pathogenesis of Type 1 Diabetes Mellitus	[With the increasing prevalence of pre-existing type 1 and type 2 diabetes in pregnancy and their associated complications, understanding the pathogenesis of these conditions is crucial.]
37121117	The association between low baroreflex sensitivity and early cardiovascular autonomic neuropathy in type 1 diabetes mellitus	[StringElement('Low baroreflex sensitivity is an indicator of early cardiovascular autonomic neuropathy. We explored the association between low baroreflex sensitivity and early cardiovascular autonomic neuropathy in type 1 diabetes mellitus.)]
37120620	Identification of Latent Autoimmune Diabetes in Adults (LADA) Using Clinical and Metabolic Features	[Latent autoimmune diabetes in adults (LADA) has clinical and metabolic features of type 1 and type 2 diabetes. LADA does not fit neatly into either category.]
37116019	A Case of Autoimmune Polyendocrine Syndrome (APS) Type II (Schmidt's Syndrome)	[Autoimmune polyendocrine syndrome (APS) type II (Schmidt's syndrome) is defined by the coexistence of autoimmune thyroid disease and other autoimmune disorders.]
37101033	Clinical Practice Guidelines for the Management of Gestational Diabetes Mellitus	[In 1989 the St. Vincent Declaration aimed to achieve comparable pregnancy outcomes in women with diabetes and those without.]
37101030	Diagnosis and Classification of Type 2 Diabetes Mellitus	[This guideline summarizes diagnosis of type 2 diabetes, including accompanying autoimmune disorders, insulin therapy, and management.]
37062046	Budget Impact Analysis of FreeStyle Libre Flash Continuous Glucose Monitoring System	[To estimate the budget impact of the potential coverage of FreeStyle Libre Flash Continuous Glucose Monitoring System.]
37046364	Use of Insulin Pumps in Type 1 Diabetes Mellitus	[The use of continuous subcutaneous insulin infusion (CSII) via insulin pumps is today considered standard of care for type 1 diabetes mellitus.]
37043824	Evaluation of the Performance of a Real-time Continuous Glucose Monitor (CGM) in Individuals with Diabetes Mellitus	[To evaluate the performance of a real-time continuous glucose monitor (CGM) in individuals with diabetes mellitus.]
37012937	A Pediatric Population-based Study of Hyperglycemia in Children with Diabetes Mellitus	[The most well-known cause of hyperglycemia is diabetes mellitus, a condition that affects the body's ability to either produce or respond to insulin.]
36945977	Low handgrip strength in older patients with type 2 diabetes who are treated with metformin	[This study aims to reveal the prevalence of low handgrip strength in older patients with type 2 diabetes who are treated with metformin.]
36920867	Novel Oncogenic Activity of Basal Insulin Fc (BIF; insulin efsitora alfa; LY3209590), a Fusion Protein Combining a Novel Single-chain Insulin Derivative with a Cytotoxic Agent	[StringElement('Basal Insulin Fc (BIF; insulin efsitora alfa; LY3209590), a fusion protein combining a novel single-chain insulin derivative with a cytotoxic agent, shows potent antitumor activity.)]
36920833	Self-Monitoring of Blood Glucose in Patients With Non-Insulin-Dependent Diabetes Mellitus.	
36916961	Adipose-tissue dysfunction and insulin resistance in type 2 diabetes mellitus	[The worldwide increase in the prevalence of diabetes mellitus (DM) has raised the demand for new therapeutic strategies.]

STAGE 2

Split Abstracts into Sentences

- Natural Language Toolkit (NLTK) library is imported and package 'punkt' is downloaded to perform sentence tokenizing using sent_tokenize function.
- Split sentences are saved in another csv file.

STAGE 2

[illegible]

STAGE 3

Text Pre-processing

- Drop empty rows
- Convert case to lowercase
- Removal of punctuations
- Removal of digit values
- Tokenization
- Removal of stopwords
- Stemming
- Lemmatization

STAGE 3

sentences

diabet one common phenotyp wolfram syndrom owe presenc variant wf gene often misdiagnos type diabet

aim explor preval wfsrelat diabet wfsdm clinic characterist chine popul earlyonset type diabet eod

sequenc exon wf gene patient eod age diagnosisuâ%âuxayear rare variant

pathogen defin accord standard guidelin american colleg medic genet genom

identifi rare variant predict deleteri patient

fast ngml postprandi cpeptid level ngml patient wf variat lower patient without wf variat respect ngml

six patient carri pathogen like pathogen variant met diagnost criterion wfsdm accord latest guidelin typic phenotyp wolfram syndrom seldom observ

diagnos earlier age usual present absenc obes impair beta cell function need insulin treatment

wfsdm usual mistakenli diagnos type diabet genet test help individu treatment

maturityonset diabet young modi group hereditari monogenet form diabet

modi account person diabet often undiagnos misdiagnos type diabet type diabet gestat diabet

diagnos modi essenti optim treatment outsid pregnanc depend modi type

review focus outcom treatment three common type modi pregnanc

diabet metabol disord glucos homeostasi Î² cell destruct occur silent detect mainli symptom appear

last year emerg great interest develop marker capabl detect pancreat Î² cell death focus improv earli diagnosi get better treatment respons mainli type diabet

type diabet would also benefit earli detect Î² cell death

differenti methyl circul dna studi minim invas biomark cell death

aim explor whether unmethylatedmethyl ratio insulin amylin gene might good biomark Î² cell death differ type diabet

lower index â†ct indic higher rate Î²cell death

plasma sampl subject without diabet pregnant woman pregnant gestat diabetesxagdm type diabet type diabet analyz

qpcr reaction specif primer methyl unmethyl fragment insulin amylin gene carri

pregnant woman gdm non gdm show higher Î²cell death marker â†insuuuÂ±u â†amylinuuuÂ±u wherea td present lower rate â†insuuuÂ±u â†amylinuuuÂ±u compar healthi subject

insulin methyl index associ newborn birth weight ruu puu insulin resist ruu puu gdm group

higher rate Î²cell death observ pregnant woman independ metabol statu

STAGE 4

Text Labelling

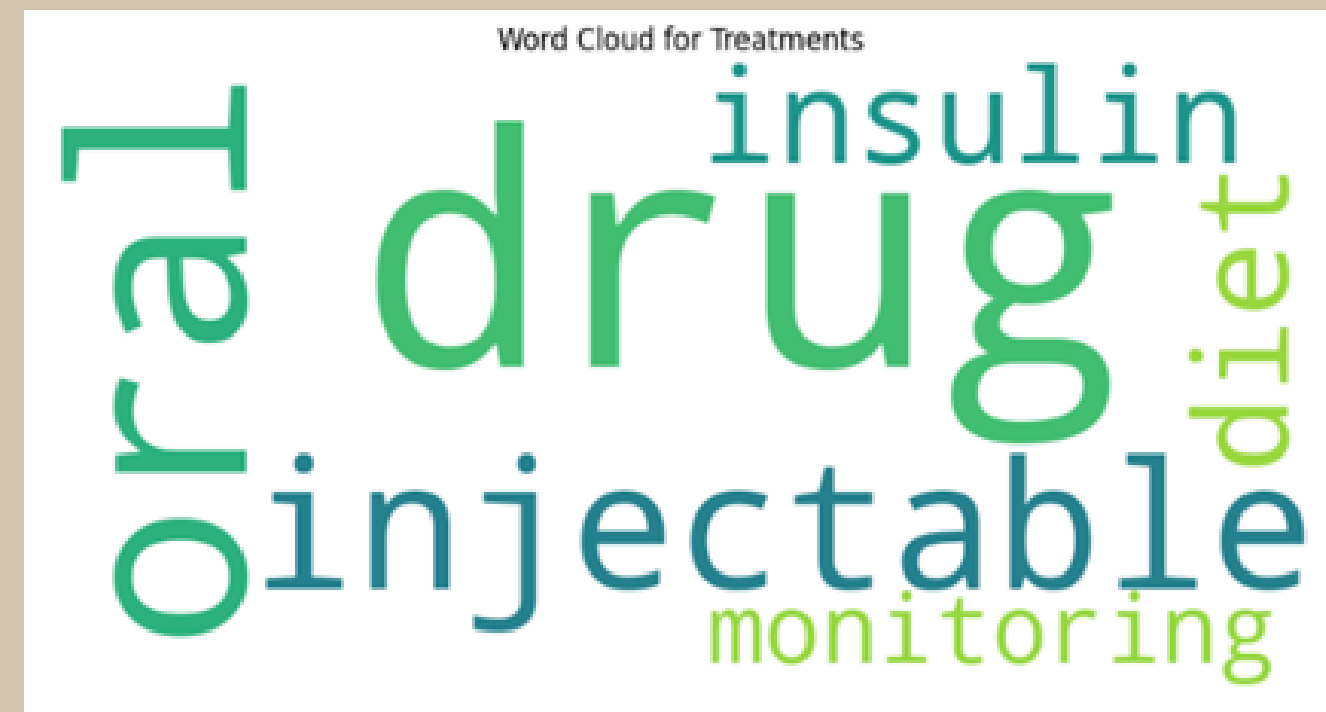
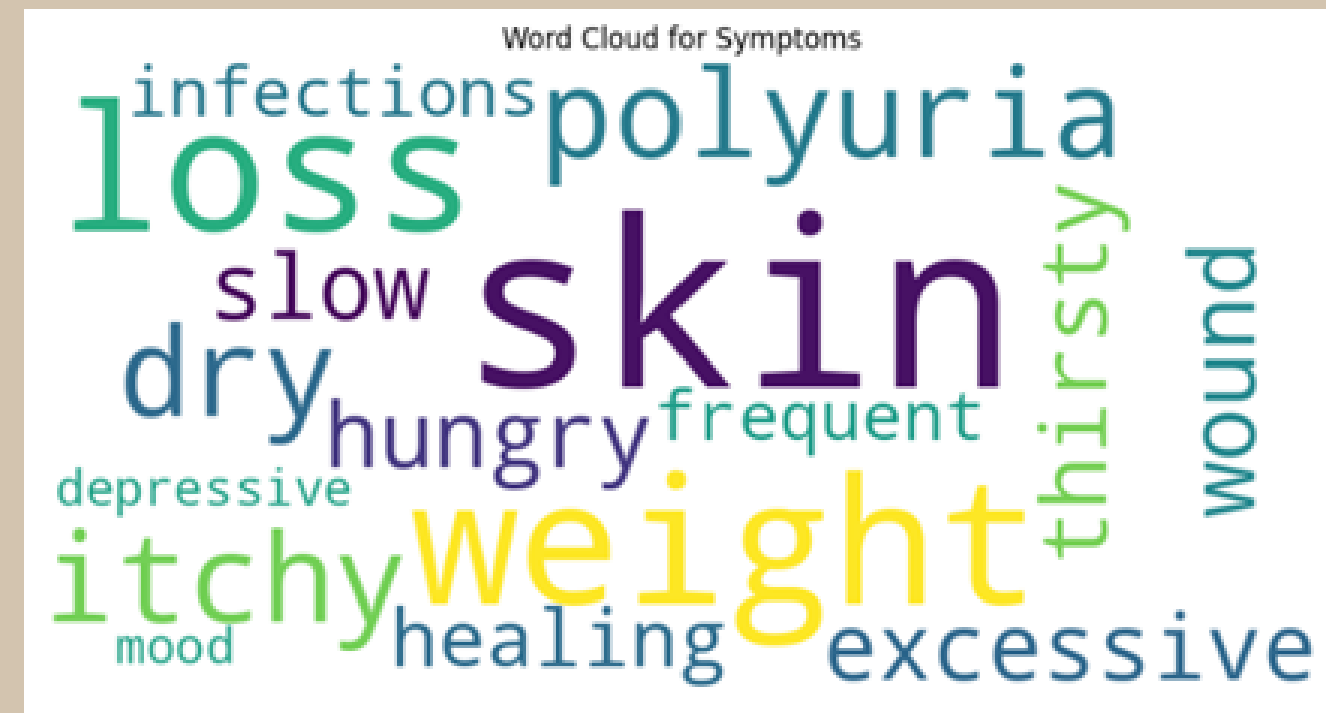
- Define Symptom and Treatment keywords for T1DM and T2DM:
 - Symptom: weight loss, polyuria, dry skin, itchy skin, excessive hungry, thirsty, blurred vision, slow wound healing, frequent infections, and depressive mood
 - Treatment: oral drug, injectable drug, insulin, and diet monitoring
- Pretrained Word2Vec: Identify synonyms of each keyword
- The sentence will be labelled as '1' if the symptoms or treatments keyword exist in the sentence; else, it will be marked as '0'

STAGE 4

[illegible]

STAGE 5

Exploratory Data Analysis





CHAPTER 5

CONCLUSION

AND

RECOMMENDATIONS

Research Outcomes

Objective 1:

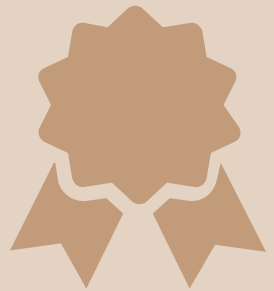


To identify the related features that are relevant to Diabetes Mellitus (DM) symptoms and treatment in multiple documents.

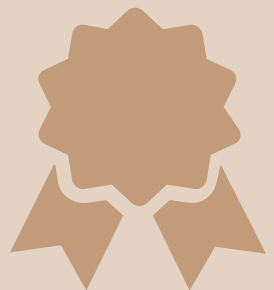
- **The textual data of DM documents was successfully scrapped from PubMed websites by using specific search terms.**
- **Each of the extracted abstracts from multiple documents were split into a set of single sentences and cleaned using text pre-processing techniques.**
- **The cleaned dataset is then labelled using the defined keywords and their similar words that obtained from the pretrained Word2Vec model.**

Achievements

Complete Data Preparation:



300 journal articles that related to DM symptoms and treatments were collected from PubMed websites.



The abstracts in the dataset were split into sentences and the cleaned dataset is obtained after text cleaning.



The sentences in cleaned dataset were labelled according to the matching of synonyms of predefined keywords for T1DM and T2DM symptoms and their treatments with the words in the sentences based on their similarity.

Future Works

In PSM II:



Identify the solutions for labelling the symptoms and treatments more precisely since the tokenization method separate the words, which cause certain words that are not actually the symptoms and treatments will also be labelled.



Develop text classification model for the collection of split and labelled sentences of Diabetes Mellitus (DM) that involves T1DM and T2DM symptoms and treatments using Support Vector Machine (SVM).



Evaluate the performance of the trained model with SVM algorithms to identify Diabetes Mellitus (DM) symptoms and treatments.



THANK YOU

End For PSM1