

**A Computational Approach for Medical Diagnosis Via the  
Decision Tree-Based Machine Learning Algorithms: A Case  
Study of Diagnosis Breast Cancer.**



**A THESIS SUBMITTED IN (PARTIAL) FULFILLMENT OF THE  
REQUIREMENT FOR THE DEGREE OF MASTER OF SCIENCE  
APPLIED MATHEMATICS  
DEPARTMENT OF MATHEMATICS, FACULTY OF SCIENCE  
KING MONGKUT'S INSTITUTE OF TECHNOLOGY LADKRABANG  
YEAR 2019  
KMITL-2019-SC-M-001-002**

การคำนวณเพื่อการวินิจฉัยทางการแพทย์บนพื้นฐานของขั้นตอน  
วิธีการเรียนรู้ของเครื่องด้วยวิธีต้นไม้การตัดสินใจ: กรณีศึกษาการ  
วินิจฉัยมะเร็งเต้านม

A Computational Approach for Medical Diagnosis Via  
the Decision Tree-Based Machine Learning Algorithms:  
A Case Study of Diagnosis Breast Cancer.



หัวข้อวิทยานิพนธ์นี้เป็นส่วนหนึ่งของการศึกษาตามหลักสูตร  
ปริญญาวิทยาศาสตรมหาบัณฑิต สาขาวิชาคณิตศาสตร์ประยุกต์  
ภาควิชาคณิตศาสตร์ คณะวิทยาศาสตร์  
สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง  
พ.ศ.2562

KMITL-2019-SC-M-001-002

This material is reserved for educational use only, not allowed for commercial use.

Forbidden to modify the content, and cite the document when use.

**A Computational Approach for Medical Diagnosis Via  
the Decision Tree-Based Machine Learning Algorithms:  
A Case Study of Diagnosis Breast Cancer.**



**A THESIS SUBMITTED IN (PARTIAL) FULFILLMENT OF THE  
REQUIREMENT FOR THE DEGREE OF MASTER OF SCIENCE  
APPLIED MATHEMATICS  
DEPARTMENT OF MATHEMATICS, FACULTY OF SCIENCE  
KING MONGKUT'S INSTITUTE OF TECHNOLOGY LADKRABANG  
YEAR 2019**

**KMITL-2019-SC-M-001-002**

This material is reserved for educational use only, not allowed for commercial use.

Forbidden to modify the content, and cite the document when use.



**COPYRIGHT 2019**

**FACULTY OF SCIENCE**

**KING MONGKUT'S INSTITUTE OF TECHNOLOGY LADKRABANG**

This material is reserved for educational use only, not allowed for commercial use.

Forbidden to modify the content, and cite the document when use.

หัวข้อวิทยานิพนธ์	การคำนวณเพื่อการวินิจฉัยทางการแพทย์บนพื้นฐานของขั้นตอนวิธีการเรียนรู้ของเครื่องด้วยวิธีต้นไม้การตัดสินใจ: กรณีศึกษาการวินิจฉัยมะเร็งเต้านม
ชื่อนักศึกษา	นางสาว จิราภรณ์ เจริญพงษ์
รหัสประจำตัว	60605091
ปริญญา	วิทยาศาสตรมหาบัณฑิต สาขาวิชาคณิตศาสตร์ประยุกต์
ภาควิชา	คณิตศาสตร์
คณะ	วิทยาศาสตร์
มหาวิทยาลัย	สถาบันเทคโนโลยีพระจอมเกล้าเจ้าคุณทหารลาดกระบัง (สจล.)
พ.ศ.	2562
อาจารย์ที่ปรึกษาวิทยานิพนธ์	ผศ.ดร.บุษยมาส พิมพ์พรรณชาติ
อาจารย์ที่ปรึกษาวิทยานิพนธ์ร่วม	รศ.ดร.วรรณพงษ์ เตรียมโพธิ์

### บทคัดย่อ

งานวิจัยนี้มีวัตถุประสงค์คือการพัฒนาแบบจำลองต้นไม้การตัดสินใจเพื่อวินิจฉัยมะเร็งเต้านม ให้มีความถูกต้องในการทำนายมากขึ้น ข้อมูลที่ใช้ในการศึกษาคือข้อมูลของผู้ที่เข้ารับการตรวจวินิจฉัยมะเร็งเต้านมจำนวน 699 คน และปัจจัยในการวินิจฉัย 9 ปัจจัย จาก UCI Machine Learning DataSet ผลการศึกษาพบว่าเมื่อนำข้อมูลมาสร้างแบบจำลองโดยใช้ต้นไม้ตัดสินใจมีความถูกต้องในการทำนาย 95.5 % โดยงานวิจัยนี้ใช้วิธีการแบ่งแยกร้อยละเพื่อแบ่งข้อมูลเป็นสองส่วนคือ ชุดฝึกฝนข้อมูลและชุดทดสอบข้อมูล เมื่อใช้วิธีดังกล่าวพบว่าแบบจำลองที่ให้ค่าความถูกต้องมากที่สุดคือ 95.4% จากนั้นเราจึงนำแบบจำลองที่ใช้ชุดฝึกฝนข้อมูล 95% มาทดสอบด้วยข้อมูล 100% พบว่าแบบจำลองดังกล่าวให้ค่าความถูกต้องเป็น 98.86% ซึ่งค่าความถูกต้องเพิ่มขึ้นจากเดิม 3.36% ถ้าหากวิธีต้นไม้การตัดสินใจมีความน่าเชื่อถือ และถ้าผลลัพธ์ของมันเป็นที่ยอมรับได้ก็อาจทำให้วิธีการนี้สามารถสร้างประโยชน์ในทางการแพทย์ได้ ผลลัพธ์ของงานวิจัยนี้พบว่าการตรวจเพียง 8 ปัจจัยก็เพียงพอต่อการวินิจฉัยมะเร็งเต้านม โดยประโยชน์ของผลลัพธ์นี้อาจช่วยในการลดงบประมาณของโรงพยาบาลในการตัดปัจจัยที่ไม่สำคัญออกไป

**คำสำคัญ :** การทำเหมืองข้อมูล การแบ่งร้อยละ การวินิจฉัย ต้นไม้การตัดสินใจ มะเร็งเต้านม

<b>Thesis Title</b>	A Computational Approach for Medical Diagnosis Via the Decision Tree-Based Machine Learning Algorithms: A Case Study of Diagnosis Breast Cancer.
<b>Student Name</b>	Miss Jiraphorn Charoenpong
<b>Student ID</b>	60605091
<b>Degree</b>	Master of Science (Applied Mathematics)
<b>Department</b>	Mathematics
<b>Faculty</b>	Science
<b>University</b>	King Mongkut's Institute of Technology Ladkrabang (KMITL)
<b>Year</b>	2019
<b>Thesis Advisor</b>	Asst. Prof. Dr. Busayamas Pimpunchat
<b>Thesis Co-Advisor</b>	Assoc. Prof. Dr. Wannapong Triampo

### Abstract

This research aims to develop the model for more accurate diagnosis based on the decision tree model and find the performance of the decision tree model that gives the best accuracy. The objective of this research is to develop the decision tree model for diagnosing breast cancer, by using data from the UCI Machine Learning DataSet from Wisconsin Repository. We found that the data 699 instances to create a decision tree without the Percentage method has the accuracy to 95.5%. We use the percentage split to develop the model. It found that the model given the highest accuracy is divided the dataset as 95% or 664 instances for creating the model and 10% or 35 instances for testing the model, which this model has 95.4% of the accuracy. Then, we brought the 699 instances to test decision trees model and provide the highest accuracy with a percentage split. The results found 98.86% of the accuracy, it increased to 3.36%. Our result, if this algorithm is more acceptable and reliable. The method can also be used in medicine and we hope to be able to help to support the decision of the physician and to help patients get faster treatment. And we hope that it can help save the budget in diagnosis because unnecessary factor has cut from our study. In this research, we found the factors no need to check for all the 9 factors but use only 8 factors were able to diagnose the disease as well, which can save budget and reduce the time to diagnosis.

**Keywords :** Breast Cancer, Data mining, Decision tree, Diagnosis, Percentage Split.

Forbidden to modify the content, and cite this document when use.

## Acknowledgements

I would like to acknowledge and thank the following important people we have supported me. Firstly, I would like to express my gratitude to my advisors, Asst. Prof. Dr. Busayamas Pimpunchat and Assoc. Prof. Dr. Wannapong Triampo for unwavering support, guidance and insight throughout this thesis. Finally, I would like to thank all my close-friends and family. They all encouraged and believed in me. They have all help me to focus on what has been a hugely rewarding and enriching process.

Jiraphorn Charoenpong



# Table of Contents

	Page
Abstract in Thai .....	i
Acknowledgements .....	iii
Table of Contents .....	iv
List of Tables.....	v
List of Figures.....	vi
<b>Chapter 1. Introduction.....</b>	<b>1</b>
1.1 Research Motivation .....	1
1.2 Objectives of the study .....	1
1.3 Scopes of the study.....	2
1.4 Research methodology.....	2
1.5 Benefits of the study .....	2
<b>Chapter 2. Methodologies and Literature reviews.....</b>	<b>4</b>
2.1 Definitions .....	4
2.1.1 Decision Tree .....	4
2.1.2 Percentage Split (Holdout Method).....	10
2.1.3 Measuring the Performance of a Classifier.....	10
2.1.4 Performance Measures .....	11
2.2 Related Research .....	12
<b>Chapter 3. Research Methodology.....</b>	<b>15</b>
3.1 Algorithms.....	15
3.2 The root node of the Decision tree .....	15
3.3 The branch of the Decision tree .....	16
3.4 The model of the decision tree .....	16
<b>Chapter 4. Results.....</b>	<b>18</b>
4.1 Data Description .....	18
4.2 To Split the data. ....	18
4.3 To find the root node .....	19
4.4 To find the branch of tree .....	23
4.5 Results of decision tree .....	24
<b>Chapter 5. Conclusion and Discussion.....</b>	<b>27</b>
Appendix A .....	31
Author Biography.....	39

# List of Tables

Table	Page
1.1 (Research methodology).....	2
2.1 Data for the Golf Example.....	7
2.2 True and False Positives and Negatives.....	11
2.3 Example of Confusion Matrix.....	11
4.1 Breast Cancer Wisconsin Dataset Attributes.....	19
4.2 The Result of Performance Measures.....	20
4.3 The Result of Decision Tree Model for the Diagnosis Breast Cancer.....	26
4.4 Confusion Matrix of Decision Tree for Diagnosis of Breast Cancer by using 35 instances.....	26
4.5 Confusion Matrix of Decision Tree for Diagnosis of Breast Cancer by using 699 instances.....	26



# List of Figures

Figure	Page
2.1 Basic Structure of Decision tree.....	4
2.2 Decision tree model that the root is outlook. ....	9
2.3 Decision Tree for the Golf Example. ....	10
3.1 Flow Chart for Diagnosis Breast Cancer by Using the Decision Tree to Building Model with Percentage Split Method. ....	15
3.2 The accuracy of training Dataset (5% to 20%) ....	16
3.3 The accuracy of training Dataset (25% to 40%) ....	17
3.4 The accuracy of training Dataset (45% to 60%) ....	17
3.5 The accuracy of training Dataset (65% to 80%) ....	17
3.6 The accuracy of training Dataset (85% to 95%) ....	17
4.1 Number of Instances in Attribute Class of Dataset. ....	18
4.2 The Accuracy of 5% to 95% of the Training Dataset.....	19
4.3 Preparing the data for find the root node.....	20
4.4 Putting the data on table for calculate the root node.....	21
4.5 Calculating the data for find the root node.....	21
4.6 Preparing the data of Clump Thick for calculate root node .....	22
4.7 The Calculation of branch.....	22
4.8 The Calculation of branch.....	24
4.9 Decision tree model for testing data set by using percentage split 95:5. ....	25

# Chapter 1

## Introduction

In this research, we will discuss the research motivation by adapting some computations of the thyroid mathematical model, the objective of the study, the scope of the study and the benefits of the study.

### 1.1 Research Motivation

Breast cancer is the top cancer in the developing world (Global Health Estimates, WHO 2015). For women, the 3 most commonly diagnosed cancers are breast, lung, bronchus, and colorectum, which collectively represent one-half of all cases; breast cancer alone is expected to account for 30% of all new cancer diagnoses in women [1].

Data mining is the process of discovering interesting patterns and knowledge from large amounts of data [2]. The decision tree is a technique in data mining, which is a method of data classification. By using this technique to diagnose breast cancer and validate the decision tree model to be able to use it. In the past several years, a decision tree was applied to diagnostics, such as Shouman M. et al. using J4.8 Algorithm Decision Tree to the diagnosis of heart disease by using other discretization techniques [3]. Sumbaly R. et al. using decision tree for diagnosis of Breast Cancer by using the J4.8 algorithm in WEKA, the results of this study are the error rates and accuracy as 5.436%, 94.564% , respectively [5]. There is also a comparison of performance by dividing the data into two parts, as a training and a testing dataset, which contain several based such as cross-validation and percentage split. Muntham D. et al. was applied the Diagnosis of the Respiratory System by using the Decision Tree algorithm i.e. ID3, C4.5, and CART and the validity of the decision tree with cross-validation and percentage split method [6]. So, we choose this method in the research.

The aim of this paper is to present the decision tree model for breast cancer. By using the decision tree method to create the model. The advantage of this method is suitable for large data. To develop a model more accuracy using percentage split. We are using the measure which included accuracy, sensitivity, and specificity for investigating the performance of the model and choose the model that gives the best performance.

### 1.2 Objectives of the study

- 1) To studying Decision Tree to diagnose breast cancer.

This material is reserved for educational use only, not allowed for commercial use.

Forbidden to modify the content, and cite the document when use.

- 2) To develop a model more accuracy by using the percentage split method.

### 1.3 Scopes of the study

- 1) To study Decision Tree and percentage split method.
- 2) To study the diagnosis of breast cancer by using the DataSet UCI from Wisconsin.
- 3) To compare the accuracy of the result in research with accuracy in other research.

### 1.4 Research methodology

- 1) Problem Identification
- 2) Studying on Decision Tree and percentage split method.
- 3) Studying on breast cancer data.
- 4) Researching about the Decision Tree for diagnosis.
- 5) Applying the Decision Tree with the breast cancer DataSet.
- 6) To analyze the result from the Decision Tree model.
- 7) Modifying the accuracy using percentage split.
- 8) Summarizing the results.

Table 1.1: (Research methodology)

Activity	Time frame				
	2018			2019	
	Apr. - Jun.	Jul. - Sep.	Oct. - Dec.	Jan. - Mar.	Apr. - Aug.
Step 1					
Step 2					
Step 3					
Step 4					
Step 5					
Step 6					
Step 7					

### 1.5 Benefits of the study

- 1) The result of this study be useful to medical for decision support of a physician.

2) To help patients get faster diagnosis of the treatment

Forbidden to modify the content, and cite the document when use.

- 3) The result of this study could help to save the budget due to the parameter of diagnosis reduce.



## Chapter 2

# Methodologies and Literature reviews

### 2.1 Definitions

#### 2.1.1 Decision Tree

Decision Tree is a tree similar hierarchical structure that compose of branches and three types of root node, internal node and leaf node respectively that correspond to the sequence of decision rules [8]. Figure 2.1 show the component of decision tree. There are many types of Decision Trees. The dissimilarity between them is the mathematical model that is used in choosing the splitting attribute in extracting the Decision Tree rules. The research tests the three most generally used types: Information Gain, Gini Index, and Gain Ratio. Different decision tree algorithms are ready to use to classify the data, such as Gini index in CART, information gain in ID3 and gain ratio C4.5 [9].

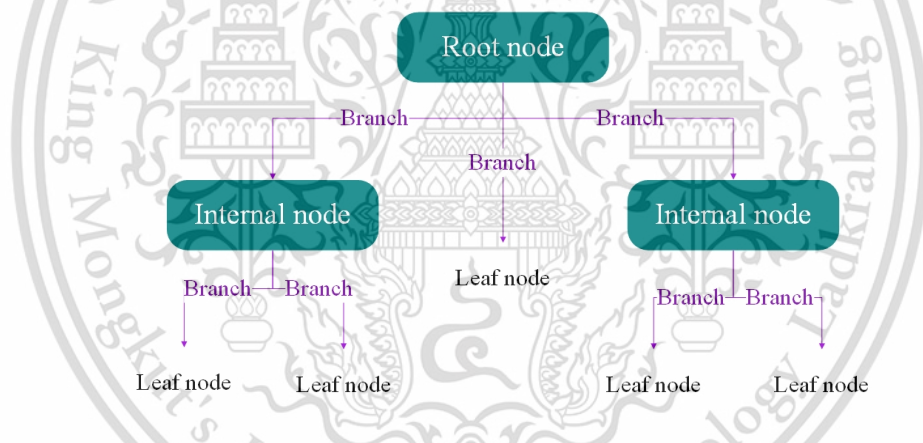


Figure 2.1: Basic Structure of Decision tree

Let the data partition  $D$  be a training set of class-labeled tuples, Suppose the class label attribute has  $m$  distinct values defining  $m$  distinct classes  $C_i$  for  $i = 1, 2, \dots, m$ . Let  $C_{i,D}$  be the set of tuples of class  $C_i$  in  $D$ . Let  $D$  and  $C_i, D$  denotes the number of tuples in  $|D|$  and  $|C_{i,D}|$ , respectively. The ID3 technique to build a decision tree is based on information theory and attempts to minimize the expected number of comparisons. Entropy is used to measure the amount of uncertainty in the data set. Defined by the formula (2.1) [3, 4, 10].

$$Entropy(D) = - \sum_{i=1}^n p_i \log_2 p_i \quad (2.1)$$

where  $p_i$  is the probability that an arbitrary tuple in  $D$  belongs to class  $C_i$ .

These conditions are:

1.  $S(p_1, p_2, \dots, p_n)$  is a continuous function.

Forbidden to modify the content, and cite the document when use.

2.  $f(n) = S(\frac{1}{n}, \frac{1}{n}, \dots, \frac{1}{n})$  is a monotonically increasing function of  $n$ .

3. Composition law for compound experiments:  $S(AB) = S(A) + \sum_{k=1}^m p_k S(B|A)$

Consider an experiment in which we randomly pick 1 object out of  $N$  objects. The probability of picking any object is  $\frac{1}{N}$ . The uncertainty of this experiment is

$$f(N) = S(\frac{1}{N}, \frac{1}{N}, \dots, \frac{1}{N}) \quad (2.2)$$

Discrimination the  $N$  objects into  $m$  groups. Each group  $k$  contains  $n_k$  objects and  $k$  runs from 1 to  $m$ ,  $\sum_{k=1}^m n_k = N$ . In the first step, resampling one of the  $m$  groups, the probability of picking group  $k$  is

$$p_k = \frac{n_k}{N} \quad (2.3)$$

Suppose group  $k$  is selected in the first step, then the probability of selecting one object in the second step is

$$\frac{1}{n_k} \quad (2.4)$$

The expected value for the uncertainty in the second step is

$$\sum_{k=1}^m p_k f(n_k) \quad (2.5)$$

Hence

$$f(N) = S(p_1, \dots, p_m) + \sum_{k=1}^m p_k f(n_k) \quad (2.6)$$

A special case of  $n_1 = n_2 = n_3 = \dots = n$ ,  $p_k = \frac{1}{m}$  for all  $k$ . Every group has  $n$  objects,  $n \cdot m = N$

$$\begin{aligned} \sum_{k=1}^m p_k f(n_k) &= p_1 f(n_1) + p_2 f(n_2) + \dots + p_m f(n_m) \\ &= p_1 f(n) + p_2 f(n) + \dots + p_m f(n) \\ &= f(n) \cdot (p_1 + p_2 + \dots + p_m) \\ &= f(n) \cdot m \left(\frac{1}{m}\right) = f(n) \end{aligned} \quad (2.7)$$

Then,

$$f(N) = S\left(\frac{1}{m}, \dots, \frac{1}{m}\right) + f(n) \quad (2.8)$$

$$f(mn) = f(m) + f(n) \quad (2.9)$$

Thus

$$f(m) = \log(m) \quad (2.10)$$

In general case

$$\begin{aligned}
 \log N &= S(p_1, \dots, p_m) + \sum_{k=1}^m p_k \log(n_k) \\
 S(p_1, \dots, p_m) &= \log N - \sum_{k=1}^m p_k \log(n_k) \\
 S(p_1, \dots, p_m) &= - \left( - \sum_{k=1}^m p_k \log N + \sum_{k=1}^m p_k \log(n_k) \right) \\
 &= - \sum_{k=1}^m p_k \log \frac{n_k}{N} \\
 &= - \sum_{k=1}^m p_k \log p_k
 \end{aligned} \tag{2.11}$$

However, it is quite likely that the partitions will be impure. Since a partition may contain a collection of tuples from different classes rather than from a single class, so must be computed  $Average(Entropy(D))$  which is the amount of data that is used to divide the data set  $D$  is a subset. This amount is measured by (2.12)[3, 4, 10]

$$Average(Entropy(D)) = - \sum_{j=1}^n \frac{|D_j|}{|D|} \times Entropy(D) \tag{2.12}$$

The term  $\frac{|D_j|}{|D|}$  acts as the weight of the  $j$ th partition.  $Average(Entropy(D))$  is the expected information required to classify a tuple from  $D$  based on the partitioning by  $A$ . The smaller the expected information required, the greater the purity of the partitions.

Information gain is defined as the difference between the original information requirement and the new requirement. After that is (2.13)[2, 3, 4, 10],

$$InformationGain(D) = Entropy(D) - Average(Entropy(D)) \tag{2.13}$$

It tells us how much would be gained by branching on  $A$ . It is the expected reduction in the information requirement caused by knowing the value of  $A$ . The attribute  $A$  with the highest Entropy, Information Gain is chosen as the splitting attribute at node  $N$ .

The C4.5 algorithm [9], a successor of ID3, uses an extension to Entropy known as a gain ratio, which attempts to overcome this bias. It applies a kind of normalization to Entropy using a “split information” value defined analogously with (2.11) as (2.14) [2, 3, 4, 10]

$$SplitInfo_A(D) = - \sum_{j=1}^n \frac{|D_j|}{|D|} \times \log_2 \left( \frac{|D_j|}{|D|} \right) \tag{2.14}$$

This value represents the potential information generated by splitting the training data set  $D$  into  $v$  partitions, corresponding to the  $v$  outcomes of a test on attribute  $A$ .

It differs from information gain, which measures the information with respect to the classification that is acquired based on the same partitioning. The gain ratio is

defined as (2.5) [4, 7]. The attribute with the maximum gain ratio is selected as the splitting attribute.

$$GainRatio(D) = \frac{InformationGain(A)}{SplitInfo_A(D)} \quad (2.15)$$

### Decision Tree: Example

In this part, we're giving an example using Decision Tree for decides golf whether or not to play each day on the basis of the weather. The data consists of 4 factors: outlooks (weather: sunny, overcast, rainy), Temp. (Temperature), Humidity and Windy. Table 2.1 shows data of 14 days of weather observation related to decision golf play. [4]

Table 2.1: Data for the Golf Example.

outlook	Temp.(F)	Humidity(%)	Windy	Class
sunny	75	70	true	play
sunny	80	90	true	don't play
sunny	85	85	false	don't play
sunny	72	95	false	don't play
sunny	69	70	false	play
overcast	72	90	true	play
overcast	83	78	false	play
overcast	64	65	true	play
overcast	81	75	false	play
rain	71	80	true	don't play
rain	65	70	true	don't play
rain	75	80	false	play
rain	68	80	false	play
rain	70	96	false	play

The creation of the decision tree model will select the attributes associated with the class as the root node by using Information Gain (IG). This value is calculated from the following equation:

$$\begin{aligned} Entropy(parent) &= -p(play) \times \log_2(play) - p(don'tplay) \times \log_2(don'tplay) \\ &= \frac{-9}{14} \times \log_2 \frac{9}{14} - \frac{5}{14} \times \log_2 \left( \frac{5}{14} \right) \\ &= 0.97 \end{aligned}$$

We calculate the entropy for the Attributes Outlook:

$$\begin{aligned}
 \text{Entropy}(\text{Outlook} = \text{sunny}) &= -p(\text{play}) \times \log_2(\text{play}) - p(\text{don't play}) \times \log_2(\text{don't play}) \\
 &= \frac{-2}{5} \times \log_2 \frac{2}{5} - \frac{3}{5} \times \log_2 \left( \frac{3}{5} \right) \\
 &= 0.97
 \end{aligned}$$

$$\begin{aligned}
 \text{Entropy}(\text{Outlook} = \text{rain}) &= -p(\text{play}) \times \log_2(\text{play}) - p(\text{don't play}) \times \log_2(\text{don't play}) \\
 &= \frac{-3}{5} \times \log_2 \frac{3}{5} - \frac{-2}{5} \times \log_2 \left( \frac{2}{5} \right) \\
 &= 0.97
 \end{aligned}$$

$$\begin{aligned}
 \text{Average}(\text{Entropy}(\text{outlook})) &= p(\text{outlook} = \text{sunny}) \times \text{Entropy}(\text{outlook} = \text{sunny}) \\
 &\quad + p(\text{outlook} = \text{overcast}) \times \text{Entropy}(\text{outlook} = \text{overcast}) \\
 &\quad + p(\text{outlook} = \text{rain}) \times \text{Entropy}(\text{outlook} = \text{rain}) \\
 &= \left( \frac{5}{14} \times 0.97 \right) + \left( \frac{4}{14} \times 0 \right) + \left( \frac{5}{14} \times 0.97 \right) \\
 &= 0.68
 \end{aligned}$$

$$\begin{aligned}
 \text{InformationGain}(\text{outlook}) &= \text{Entropy}(\text{parent}) - \text{Average}(\text{Entropy}(\text{outlook})) \\
 &= 0.94 - 0.68 \\
 &= 0.24
 \end{aligned}$$

$$\begin{aligned}
 \text{InformationGain}(\text{Temperature}) &= \text{Entropy}(\text{parent}) - \text{Average}(\text{Entropy}(\text{Temperature})) \\
 &= 0.94 - 0.91 \\
 &= 0.03
 \end{aligned}$$

$$\begin{aligned}
 \text{InformationGain}(\text{Humidity}) &= \text{Entropy}(\text{parent}) - \text{Average}(\text{Entropy}(\text{Humidity})) \\
 &= 0.94 - 0.79 \\
 &= 0.15
 \end{aligned}$$

$$\begin{aligned}
 \text{InformationGain}(\text{Windy}) &= \text{Entropy}(\text{parent}) - \text{Average}(\text{Entropy}(\text{Windy})) \\
 &= 0.94 - 0.89 \\
 &= 0.05
 \end{aligned}$$

The outlook is the higher information gain, then outlook will be the root of tree. Figure 2.2 Show the decision tree model that the root is outlook. Next, we find the information gain for the rain branch.

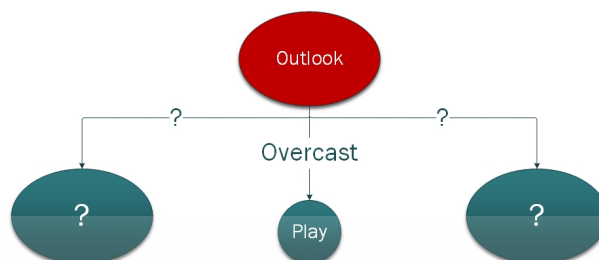


Figure 2.2: Decision tree model that the root is outlook.

$$\begin{aligned}
 Entropy(\text{parent}) &= -p(\text{play}) \times \log_2(\text{play}) - p(\text{don'tplay}) \times \log_2(\text{don'tplay}) \\
 &= \frac{-3}{5} \times \log_2 \frac{3}{5} - \frac{-2}{5} \times \log_2 \left( \frac{2}{5} \right) \\
 &= 0.97
 \end{aligned}$$

$$\begin{aligned}
 InformationGain(\text{Temperature}) &= Entropy(\text{parent}) - Average(Entropy(\text{Temperature})) \\
 &= 0.97 - 0.95 \\
 &= 0.02
 \end{aligned}$$

$$\begin{aligned}
 InformationGain(\text{Humidity}) &= Entropy(\text{parent}) - Average(Entropy(\text{Humidity})) \\
 &= 0.97 - 0.95 \\
 &= 0.02
 \end{aligned}$$

$$\begin{aligned}
 InformationGain(\text{Windy}) &= Entropy(\text{parent}) - Average(Entropy(\text{Windy})) \\
 &= 0.97 - 0 \\
 &= 0.97
 \end{aligned}$$

Windy is the higher information gain, then it will be the internal node of rain branch. we find the information gain for the sunny branch.

$$\begin{aligned}
 Entropy(\text{parent}) &= -p(\text{play}) \times \log_2(\text{play}) - p(\text{don'tplay}) \times \log_2(\text{don'tplay}) \\
 &= \frac{-2}{5} \times \log_2 \frac{2}{5} - \frac{-3}{5} \times \log_2 \left( \frac{3}{5} \right) = 0.97
 \end{aligned}$$

This material is reserved for educational use only, not allowed for commercial use.

Forbidden to modify the content, and cite the document when use.

$$\begin{aligned}
 \text{InformationGain}(\text{Temperature}) &= \text{Entropy}(\text{parent}) - \text{Average}(\text{Entropy}(\text{Temperature})) \\
 &= 0.97 - 0.4 \\
 &= 0.57
 \end{aligned}$$

$$\begin{aligned}
 \text{InformationGain}(\text{Humidity}) &= \text{Entropy}(\text{parent}) - \text{Average}(\text{Entropy}(\text{Humidity})) \\
 &= 0.97 - 0 \\
 &= 0.97
 \end{aligned}$$

Humidity is the higher information gain, then it will be the internal node of sunny branch. Decision tree model of golf playing decisions as shown in Figure 2.3. From Figure 2.3, the model can be described as follows: if outlook is overcast or cloudy, then the golfer will decide to play golf. If outlook is sunny and Humidity more than 75 then golfers will decide not to play golf.

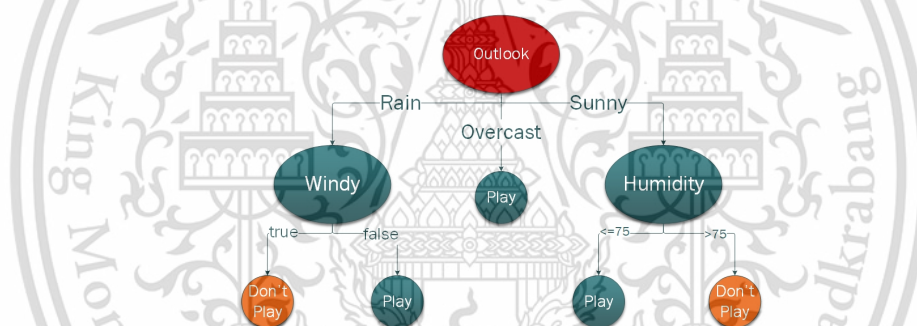


Figure 2.3: Decision Tree for the Golf Example.

### 2.1.2 Percentage Split (Holdout Method)

In this method, the given data are randomly partitioned into two independent sets, a training, and a test set. Typically, two-thirds of the data are allocated to the training set, and the remaining one-third is allocated for testing. The training set is used to derive the model, whose accuracy is estimated with the test set [2].

### 2.1.3 Measuring the Performance of a Classifier

Confusion Matrix

There are two classes, which call positive and negative, the confusion matrix consists of four cells, which can be labeled TP, FP, FN, TN as in Table 2.2 [2, 3, 10].

Where TP (True Positive): the number of positive instances that are classified as positive. TN (True Negative): the number of negative instances that are classified as negative. FP (False Positive): the number of negative instances that are classified as positive. FN (False Negative): the number of positive instances that are classified as negative.

Forbidden to modify the content, and cite the document when use.

**Table 2.2:** True and False Positives and Negatives.

Confusion Matrix	Predicted Class	
	Positive	Negative
Actual Class	TP	FP
	FN	TN

positive. FN (False Negative): the number of negative instances that are classified as negative.

#### 2.1.4 Performance Measures

It is the task of classifying the elements of a given set into more than one group on the basis of a classification rule. There are many metrics that can be used to measure the performance of a classifier or predictor; different fields have different preferences for specific metrics due to different goals. The accuracy of a classifier is a value that indicates the ability of a measurement. Sensitivity measures the proportion of actual positives that are correctly identified as such. Specificity measures the proportion of actual negatives that are correctly identified as such. These measures are defined as follows [2, 3, 10, 11]

$$\text{Sensitivity} = \frac{TP}{TP + FN} \quad (2.16)$$

$$\text{Specificity} = \frac{TN}{FP + TN} \quad (2.17)$$

$$\text{Accuracy} = \frac{TP + TN}{TP + FN + FP + TN} \quad (2.18)$$

To provide more insight. We give an example: Examples of breast cancer diagnoses, which have two alternatives, are malignant and benign. This data contains 699 patients consist of benign 458 and malignant 241 patients. From the prediction, the model can be evaluated using a confusion matrix, as shown in Table 2.3.

**Table 2.3:** Example of Confusion Matrix.

Confusion Matrix	Predicted Class	
	Positive	Negative
Actual Class	443(TP)	6(FP)
	15(FN)	235(TN)

From Table 2.3, it was found that 443 patients correct predicted that benign and wrong predicted of 15 patients (the patients are benign, but the model is predicted to be malignant). And, it was found that 235 patients correct predicted that malignant and

This material is reserved for educational use only, not allowed for commercial use.

Forbidden to modify the content, and cite the document when use.

wrong predicted of 6 patients (the patients are malignant, but the model is predicted to be benign). The model can be evaluated by using the accuracy value which is approximately 97

## 2.2 Related Research

Machine Learning can be applied in many things such as industry, agriculture, medical, etc. Machine Learning used for the analysis of the importance of clinical parameters and their combinations for prognosis, e.g. prediction of disease progression, extraction of medical knowledge for outcome research, therapy planning and support, and the overall patient management [12]. In the past many years, there is a lot of research to use Machine Learning Algorithm to apply in cancer. In the survival analysis for cancer. This method is used to find the survival rate of cancer patients. Application of Machine Learning to Survival will help predict survival rates. This application can be developed for survival prediction in rare pathologies and have the potential to serve as the basis for the creation of decision support tools in the future [13]. At present, the software helps us to work more conveniently and save more time. And some types of software are also available for free. In the ML, there are several types of software available for use. Weka is a tool for the initial processing of data. Advantages of Weka Firstly, Weka is an open-source. Secondly, it solves the ML problem. Finally, it runs on many platforms [14]. It also allows users to select the most appropriate algorithm for the model to obtain the most accurate value [15]. Another tool that is used frequently, the R programming. R is used in machine learning for classification, regression, and clustering. There are a variety of packages in R to make ML easier to use, such as rpart [16] is a package that helps in recursive partitioning for classification, regression and survival trees. Another important package is e1701 [17]. It can be used in support vector machines, shortest path computation, bagged clustering, and Naive Bayes classifier.

Machine Learning Applications with Cancer From survey research of machine learning application with cancer. The first matter to be studied is diagnostic. A comparison of the performance of different machine learning algorithms has been studied in many researches. This research uses public information from the UCI. This data is a breast cancer test with two classes for decision making, with malignant and benign. This information is available for 699 patients, consisting of 241 of malignant and 458 of benign. The main purpose of this research is to compare and evaluate different machine learning. By using SVM, DT (C4.5), NB and k-NN to predict the diagnosis of breast cancer and comparative efficiency [18]. The results of such research showed that SVM gives the highest accuracy (97.13%) and the lowest error rate (2%) compared with other algorithms. It can be said that the SVM algorithm is a method that provides the most

This material is reserved for educational use only, not allowed for commercial use.

Forbidden to modify the content, and cite the document when use.

accurate of this dataset for prediction diagnosis. Using ANN and LR to developing a model to predict liver cancer [19] is one of the interests by using 70% of the data to training and 30% for testing the model. The result shows that the accuracy of ANN better than the LR. In the case of survival and recurrence, there is interesting research as well. The one paper report is the prediction model of survival and recurrence for breast cancer by using machine learning [20]. The result of predict survival shows that the ANN algorithm obtained the best accuracy (83.6%) using live-one-out cross-validation as a test method. Meanwhile in the prediction of recurrence of breast cancer found that the DT algorithm has achieved the best accuracy (96.57%) by using validation like that the validation of prediction of survival. The comparison of the methods of machine learning is still many in many researches. This research uses 3 machine learning algorithms (DT (C4.5), SVM, ANN) for medical data analysis which is big data to find hidden knowledge. The data for this study is the data of breast cancer patients admitted to the Iran Center program, which has a total of 1189 patients. The results of this study showed that SVM predicted breast cancer models better than other methods with the highest correctness of 95.7% and DT with the lowest predictive value of 93.6% [21]. The next research is also a study of breast cancer. This research uses 7 machine learning algorithms consisting of NB, Trees Random Forest, 1-Nearest Neighbor, AdaBoost, SVM, RBF Network and Multilayer Perceptron. The objective of this research is to apply the seven methods of machine learning to predict the model of survival of breast cancer. The data from the study of Cancer Registry Organization of Kerman, Iran with 900 breast cancer patients, consisting of 803 living patients and 97 patients died. The results of this study showed that the Trees Random Forest predicted the best breast cancer survival model with 96% accuracy. And 1-Nearest Neighbor predicted the model to be less accurate than other methods, 91% of accuracy [22].

Shouman M. et. al. using J4.8 Algorithm Decision Tree and other discretization techniques to analyse a model of the diagnosis heart disease. The objective of this research is to show a new model that enhances the Decision Tree accuracy by apply voting. The data used in this research is heart disease data from the Cleveland Clinic Foundation. The study found that the application of Voting with the Decision Tree shows the increased accuracy [3].

Sumbaly R. et. al. using decision tree for diagnosis of Breast Cancer by using the J4.8 algorithm in WEKA. The data used in this research is Breast Cancer DataSet from the UCI Machine Learning Repository. This DataSet consist of 9 factor to diagnosis breast cancer. The results of this study are the error rates and accuracy as 5.436%, 94.564% respective. And this research also found that the diagnosis using only 5 factors are enough to diagnose breast cancer [5].

Muntham D. et. al. was applied Diagnosis of the Respiratory System by using the Decision Tree algorithm i.e. ID3, C4.5, and CART and the validity of the decision

tree with cross-validation and percentage split method. The results showed that for a patients with acute with a ratio 70:30 of the training data and the testing data, Decision Tree (C4.5) was the most effective (92.31%). For patients with acute sinusitis with ratio 70:30 of training data and testing data, Decision Tree(C4.5) was most effective (94.70%). For patients with pneumonia with ratio 50:50 of training data and testing data, the CART Algorithm was effective (94.69%) [6].

Parveen et.al studied Decision Tree to evaluate the risk of Non-Alcoholic Fatty Liver Disease in the Canadian population. The data used in this study is the risk factor of Non-Alcoholic Fatty Liver Disease from Electronic Medical Records (EMRs). The result shows that the Decision Tree method was most effective for the diagnosis of Non-Alcoholic Fatty Liver Disease risk, 93.7 of accuracy [23].

Cirkovic et.al. showed prediction model of survival and recurrence of breast cancer by using machine learning. The result of predict survival shows that the ANN algorithm obtained the best accuracy (83.6%) with using live-one-out cross-validation as a test method. Meanwhile in prediction of recurrence of breast cancer found that Decision Tree algorithm has achieved the best accuracy (96.57%) by using validation [24].

P.Hamsagayathri et.al. studied Decision Tree algorithms for Breast cancer to analyze the performance of classification. The data used in this research was study from SEER breast cancer dataset. And this research use WEKA software to breast cancer diagnosis classification. The results show that Decision Tree classifier classifies the data with 98.51% accuracy [25].

Chaurasia V. et.al used three algorithms (Naïve Bayes, Radial Basis Functions Neural Networks, J48) to develop the prediction models for breast cancer diagnosis. The data used in this research is Breast Cancer DataSet from the UCI Machine Learning Repository. This DataSet consist of 683 breast cancer cases. The results showed that Naïve Bayes is the best predictor with 97.36% accuracy. Meanwhile J48 came out the accuracy with 93.41% [26].

## Chapter 3

# Research Methodology

### 3.1 Algorithms

The algorithm to create a model for diagnosis of breast cancer by using a decision tree and choose the model that, given the best performance is shown in Figure 3.1. The steps to create a model are as follows.

- 1) Studying on Breast cancer DataSet from UCI.
- 2) Studying on classification algorithm.
- 3) Finding suitable methods for classifying with DataSet.
- 4) Studying on percentage split or holdout method
- 5) Dividing the data into 2 sets: training set and testing set
- 6) Using the training set to learning process with a decision tree to create a prediction model.
- 7) Using the testing set to test model.
- 8) Performance evaluation of breast cancer diagnosis model by using confusion matrix

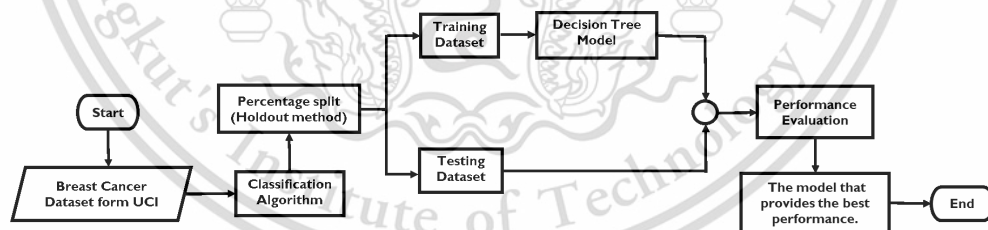


Figure 3.1: Flow Chart for Diagnosis Breast Cancer by Using the Decision Tree to Building Model with Percentage Split Method.

### 3.2 The root node of the Decision tree

Details of the model are shown in Figure 4.9. By the C5.0 algorithm, the attribute with the maximum gain ratio is selected as the splitting attribute [10]. Therefore, the attribute with the highest gain ratio is used as a root node of the tree.

### 3.3 The branch of the Decision tree

The data used to create a Decision Tree model in this paper is numeric. It consists of the following steps:

- i) Sort smallest to largest,
- ii) Find the midpoint between two different values, i.e. the midpoint between 1, 2 is 1.5, etc., then divide the data into two sections, divided into more than and less than or equal.
- iii) Calculating Entropy of two parts, the part of less than or equal of the midpoint to be on the left side of the node and the part more than to be on the right side.
- iv) Find the midpoint of all the values in the data, and then select the value to give the Gain Ratio most.

### 3.4 The model of the decision tree

To choose a decision tree model. We have divided the data into training Dataset and testing Dataset. The criterion for choosing a model is the accuracy of the model. From Figure 3.2 shows the accuracy of each model.

The figure displays four tables showing the accuracy of a training dataset for different median values (5, 10, 15, and 20). Each table has columns for 'Median', 'Accuracy', and 'Error'. The data is as follows:

Median	Accuracy	Error
5	0.95	0.05
10	0.90	0.10
15	0.85	0.15
20	0.80	0.20

Figure 3.2: The accuracy of training Dataset (5% to 20%)

From Figures 3.2 to 3.6, it shows the accuracy of dividing data into two sets. The principle of the models using the percentage split method which is divide the data into two parts and both parts must not have duplicate data.

Model (%)	Accuracy (%)	Model (%)	Accuracy (%)	Model (%)	Accuracy (%)	Model (%)	Accuracy (%)
25	93.5	30	93.75	35	94.03	40	94.04
26	93.5	31	93.75	36	94.03	41	94.04
27	94.70	32	94.47	37	94.10	42	94.20
28	92.92	33	92.03	38	92.15	43	93.03
29	93.7	34	94.43	39	94.03	44	93.04
30	93.22	35	93.22	40	93.22	45	93.22
31	94.00	36	90.17	41	90.000	46	94.100
32	92.37	37	89.57	42	90.055	47	91.786
33	93.3	38	90.96	43	90.604	48	91.466
34	93.22	39	93.22	44	93.22	49	93.22
35	94.00	40	91.13	45	90.492	50	90.040
36	91.87	41	87.11	46	87.851	51	90.506
37	94.7	42	86.60	47	87.3		
38	93.22	43	93.22				
39	94.00	44					

Figure 3.3: The accuracy of training Dataset (25% to 40%)

Model (%)	Accuracy (%)	Model (%)	Accuracy (%)	Model (%)	Accuracy (%)	Model (%)	Accuracy (%)
45	93.84	50	93.98	55	94.04	60	94.04
46	94.111	51	94.01	56	94.04	61	94.12
47	93.50	52	92.40	57	93.24	62	93.10
48	94.01	53	91.62	58	92.96	63	92.84
49	93.265	54	93.70	59	93.20	64	93.20
50	92.28	55	94.92	60	94.00	65	94.00
51	92.40	56	93.05	61	93.02	66	93.00
52	92.70	57	93.11	62	93.12	67	93.12
53	93.02	58	92.22	63	93.00	68	93.12
54	84.33	59	93.22	64	93.00	69	93.12
55	93.51	60		65		70	93.12

Figure 3.4: The accuracy of training Dataset (45% to 60%)

Model (%)	Accuracy (%)	Model (%)	Accuracy (%)	Model (%)	Accuracy (%)	Model (%)	Accuracy (%)
65	94.93	70	95.03	75	95.03	80	95.03
66	94.94	71	95.03	76	95.03	81	95.03
67	95.00	72	95.03	77	95.03	82	95.03
68	95.75	73	95.03	78	95.03	83	95.03
69	95.36	74	95.03	79	95.03	84	95.03
70	93.41	75	95.03	80	95.03	85	95.03
71	93.25	76		81		86	95.03

Figure 3.5: The accuracy of training Dataset (65% to 80%)

Model (%)	Accuracy (%)	Model (%)	Accuracy (%)	Model (%)	Accuracy (%)
85	94.98	90	95.24	95	94.82
86	94.11				
87	93.73				

Figure 3.6: The accuracy of training Dataset (85% to 95%)

From the above figures, it shows that the test data decreases, the accuracy decreases. In other words, Training data may not cover all testing data. While the training data is the same, the test data has changed. Therefore, if the data in the test increases, it may make the model more accurate. Because there are cases that can cover the model.

## Chapter 4

### Results

#### 4.1 Data Description

The data used in this study is Breast Cancer Wisconsin Data Set from UCI Machine Learning Repository [10]. This dataset having 699 instances consist of benign 458 instances and malignant 241 instances. Figure 4.1. shows the number in the attribute class of the dataset.

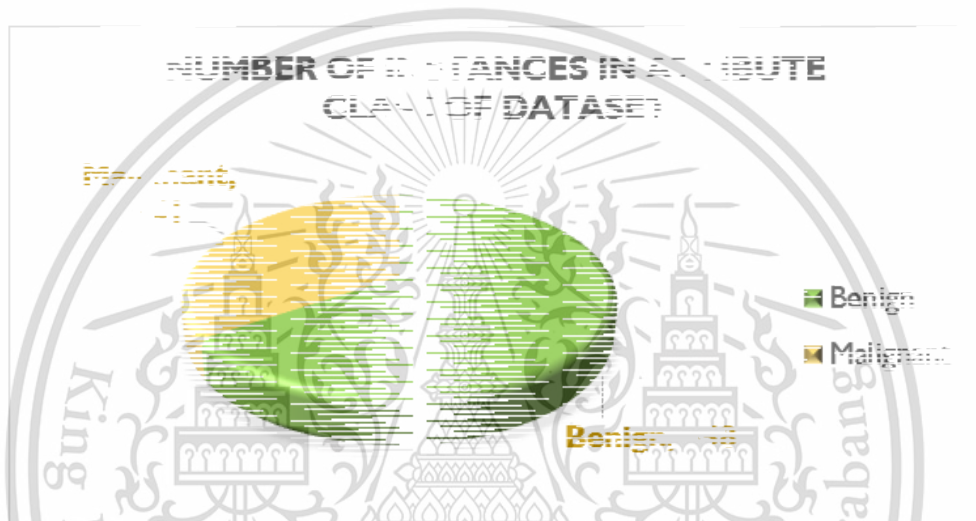


Figure 4.1: Number of Instances in Attribute Class of Dataset.

The attributes in this dataset are shown in Table 4.1. By using the attributes to diagnose, only 9 attributes because the sample code number is not important for the diagnosis of breast cancer.

#### 4.2 To Split the data.

In this study, we used the C4.5 algorithm to model building the decision tree. We bring data 699 instances to predict the diagnosed of breast cancer with a decision tree and not yet use the Percentage split method found the accuracy is 95.5%. To develop the performance of the decision tree model, we use the Percentage Split method or Holdout method [6, 23]. The Percentage Split method by determining the default ratio is 5:95 between the training dataset and the testing dataset and add each step, 5% of the training dataset and reduce each step, 5% of the testing dataset. Since the Percentage split data is divided into 2 parts, each part is selecting data using random. Therefore, to more accurate and to the more distribution of data, we randomized data 100 times then take to make a decision tree. With this technique, it brings up more

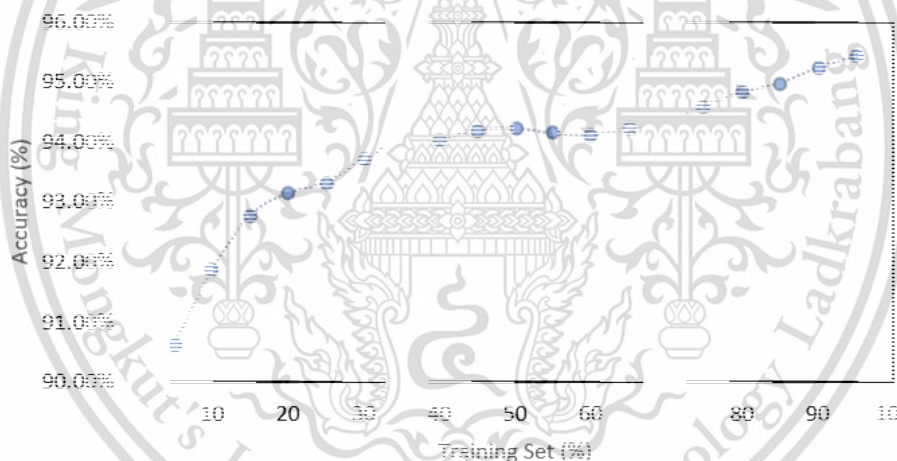
This material is reserved for educational use only, not allowed for commercial use.

Forbidden to modify the content, and cite the document when use.

**Table 4.1:** Breast Cancer Wisconsin Dataset Attributes.

No	Attribute Name	Value
1	Uniformity of Cell Shape (UoCS)	1-10
2	Uniformity of Cell Size (UoCZ)	1-10
3	Bare Nuclei (BN)	1-10
4	Bland Chromatin (BC)	1-10
5	Single Epithelial Cell Size (SECZ)	1-10
6	Normal Nucleoli (NN)	1-10
7	Marginal Adhesion (MA)	1-10
8	Clump Thickness (CT)	1-10
9	Mitoses (M)	1-10
10	Sample code number	The id number of instances
11	Class	2-Benign and 4-Malignant

accuracy, and specificity sensitivity value to find the average.

**Figure 4.2:** The Accuracy of 5% to 95% of the Training Dataset.

By data showing trends in Figure 4.2, it is seen that, agrees that accuracy tends to rise when the training dataset is increasing while the testing dataset is decreasing. It was found that a ratio 95:5 of the training dataset and the testing dataset which provides the highest accuracy, sensitivity, and specificity are 95.4%, 95.8 and 94.6% respectively, we conclude the information shown in Table 4.2.

### 4.3 To find the root node

Details of the model are shown in Figure 4.9. By the C5.0 algorithm, the attribute with the maximum gain ratio is selected as the splitting attribute [10]. Therefore, the attribute with the highest gain ratio is used as a root node of the tree.

This material is reserved for educational use only, not allowed for commercial use.

Forbidden to modify the content, and cite the document when use.

**Table 4.2:** The Result of Performance Measures.

Training Data(%)	Testing Data(%)	Accuracy(%)	Sensitivity (%)	Specificity (%)
5	95	90.64	95.26	81.88
10	90	91.87	95.35	85.28
15	85	92.78	95.42	87.73
20	80	93.15	95.71	88.30
25	75	93.33	95.33	89.51
30	70	93.72	94.80	91.66
35	65	94.03	94.95	92.26
40	60	94.04	94.98	92.25
45	55	94.21	95.06	92.56
50	50	94.23	94.58	93.60
55	45	94.15	94.53	93.41
60	40	94.12	94.46	93.48
65	35	94.23	94.67	93.44
70	30	94.35	94.34	94.43
75	25	94.61	94.87	94.11
80	20	94.83	94.98	94.56
85	15	94.98	95.48	94.06
90	10	95.24	95.57	94.48
95	5	95.43	96.80	94.64

It consists of the following steps:

- i) Preparing the data as the figure 4.3

Class	Group	Form of Cell Size (Umum)	Maximal Adhesion (Single Epithelium)	Large Nuclei	Bland Chromatin	Clumps
Benign		1	1	1		1
Benign		4	5	10		1
Benign		1	1	2		1
Benign		8	1	4		1
Benign		0	3	1		1
Malignant		10	8	10		1
Benign		1	1	10		1
Benign		1	1	1		1
Benign		1	1	1		1
Benign		1	1	1		1
Benign		1	1	1		1
Benign		1	1	1		1
Benign		1	1	1		1
Malignant		1	5	3		4
Benign		1	1	3		1
Malignant		4	4	1		1
Benign		1	1	1		1
Benign		1	1	1		1
Malignant		7	6	10		2
Benign		1	1	1		1
Malignant		3	10	10		4
Malignant		5	3	7		1
Benign		1	1	1		1
Malignant		4	1	0		1
Benign		1	1	1		1
Malignant		2	4	7		1
Benign		1	1	1		1
Benign		1	1	1		1
Benign		1	1	1		1
Benign		1	1	1		1
Benign		1	1	1		1
Benign		1	1	1		1

**Figure 4.3:** Preparing the data for find the root node.

This material is reserved for educational use only, not allowed for commercial use.

Forbidden to modify the content, and cite the document when use.

ii) Taking the data to table for calculate as the figure 4.4

Figure 4.4: Putting the data on table for calculate the root node.

iii) Calculating the Entropy, Information gain, Split Information and Gain Ratio as figure 4.5

Figure 4.5: Calculating the data for find the root node.

For example, Finding the Gain Ratio of Clump Thick. Firstly, preparing the data of Uniformity of Cell Size as figure 4.6.

Next, calculating the entropy, information gain, split information and gain ratio respectively. This calculation as follows : Using the training dataset which has 629 instances detect a benign of 412 instances and malignant 217 instances by calculating Entropy of 629 instances as follows

$$Entropy(629instances) = -\frac{412}{629} \log_2\left(\frac{412}{629}\right) - \frac{217}{629} \log_2\left(\frac{217}{629}\right) = 0.929516$$

Finding the entropy of Clump Thick as follow

$$Entropy(ClumpThick) = Entropy(Level1) + Entropy(Level2) + Entropy(Level3) + Entropy(Level4) + Entropy(Level5) + Entropy(Level6) + Entropy(Level7) + Entropy(Level8) + Entropy(Level9) + Entropy(Level10) = 0.033 + 0.0188 + 0.0778 + 0.0686 + 0.1796 + 0.0475 + 0.009 + 0.0248 + 0 + 0 =$$

This material is reserved for educational use only, not allowed for commercial use.

Forbidden to modify the content, and cite the document when use.

level	Class	Uniformity of Cell Size	total	entropy
1	Maligna	342	346	
2	Maligna	32	39	
3	Maligna	29	47	
4	Maligna	8	36	
5	Maligna	0	27	
6	Maligna	2	24	
7	Maligna	1	16	
8	Maligna	1	25	
9	Maligna	1	6	
10	Maligna	0	62	
total		411	628	
entropy of CT				
Information				
Split				
Gain Ratio				

Figure 4.6: Preparing the data of Clump Thick for calculate root node

0.459150376

Then,  $InformationGain(CT) = Entropy(629ins.) - Entropy(CT) = 0.47036563$

And  $SplitInfo_A(ClumpThick) \approx 3.026918181$ .

Hence,  $GainRatio(ClumpThick) = \frac{InformationGain(ClumpThick)}{SplitInfo_A(ClumpThick)} = 0.155394176$

level	Class	Thick	entropy
	Maligna	sign tot	
		133	0.033
		43	0.0188
		83	0.0778
		56	0.0686
		77	0.1796
		16	0.0475
		1	0.009
		3	0.0248
		0	0
		0	0
total		412	
entropy		1591503	
Informat		1703654	
Split		3269181	
GainRa		1553941	

Figure 4.7: The Calculation of branch

Similarly, we can compute the Gain Ratio of other factors, in the same way, which did above. Which has the following value

Gain Ratio (Uniformity of Cell Size) = 0.293898297 , Gain Ratio (Uniformity of Cell Shape)=0.263566547, Gain Ratio(Marginal Adhesion)=0.204129124 , Gain Ratio(Bare

Nuclei)=0.20857053 , Gain Ratio(Bland Chromatin)=0.204345072 , Gain Ratio(Normal Nucleoli)=0.23180457 , Gain Ratio(Mitoses)=0.180877745 and Gain Ratio (Single Epithelial Cell Size)=0.221127018 .

Figure 4.5 shows the detail of the root node calculate. Since Gain Ratio(Uniformity of Cell Size) Given the most Gain ratio, so the root node chooses the Uniformity of Cell Size as a root.

#### 4.4 To find the branch of tree

The data used to create a Decision Tree model in this paper is numeric. It consists of the following steps:

- i) Sort smallest to largest,
- ii) Find the midpoint between two different values, i.e. the midpoint between 1, 2 is 1.5, etc., then divide the data into two sections, divided into more than and less than or equal.
- iii) Calculating Entropy of two parts, the part of less than or equal of the midpoint to be on the left side of the node and the part more than to be on the right side.
- iv) Find the midpoint of all the values in the data, and then select the value to give the Gain Ratio most.

To get a better understanding, I'll give the example for calculating the branch of a root node, which is the Uniformity of Cell Shape. The data in the Uniformity of Cell Shape range from 1 to 10. It sorts the data from the smallest to largest and then find the midpoint between two different values.

Consequently, the midpoint in the Uniformity of Cell Size has 1.5, 2.5, 3.5, 4.5, ..., 9.5 respectively. Bring up the midpoint obtained to divide the two ranges, such as  $>1.5, \leq 1.5$  etc. For creating a decision tree model, using the training dataset which has 629 instances detect a benign of 412 instances and malignant 217 instances by calculating Entropy of 629 instances as follows

$$Entropy(629instances) = -\frac{412}{629} \log_2\left(\frac{412}{629}\right) - \frac{217}{629} \log_2\left(\frac{217}{629}\right) = 0.929516$$

In the section of  $\leq 1.5$  detects benign 314 instances, malignant 2 instances and in the section of  $>1.5$  detects benign 98 instances, malignant 215 instances. To find the Gain Ratio can be calculated as follows

$$Entropy(\leq 1.5) = \frac{316}{629} \times \left(-\frac{314}{316} \log_2\left(\frac{314}{316}\right) - \frac{2}{316} \log_2\left(\frac{2}{316}\right)\right)$$

$$Entropy(> 1.5) = \frac{313}{629} \times \left(-\frac{98}{313} \log_2\left(\frac{98}{313}\right) - \frac{215}{313} \log_2\left(\frac{215}{313}\right)\right)$$

So,  $Entropy(1.5) = 0.027796 + 0.446221 = 0.474017$

Then,  $InformationGain(1.5) = Entropy(629ins.) - Entropy(1.5) = 0.455499$ .

And  $SplitInfo_A(1.5) \approx 0.999984$ .

Hence,  $GainRatio(1.5) = \frac{InformationGain(1.5)}{SplitInfo_A(1.5)} = 0.507076$

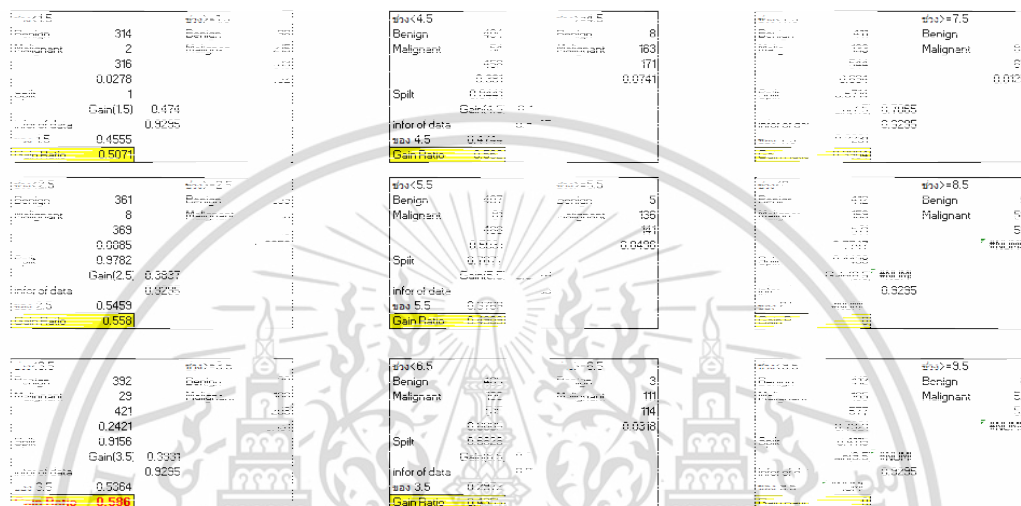


Figure 4.8: The Calculation of branch

Similarly, we can compute the Gain Ratio of other values, in the same way, which did above. Which has the following value Gain Ratio (2.5)= 0.558009 , Gain Ratio (3.5)=0.585805 ,Gain Ratio(4.5)=0.561974, Gain Ratio(5.5)=0.490921,Gain Ratio(6.5)=0.435272 ,Gain Ratio(7.5)=0.390407,Gain Ratio(8.5)=0 and Gain Ratio (9.5)=0. Figure 4.3 shows the detail of the branch of tree calculate. Since Gain Ratio(3.5)Given the most Gain ratio, so the root node chooses 3.5 as a branch.

#### 4.5 Results of decision tree

From Figure 4.9 the decision tree model can be converted to classification rules, which have 13 rules. When a biopsy was found Uniformity of Cell Size > 3.5, Uniformity of Cell Shape > 1.5 and Bland Chromatin > 1.5 have the opportunity to malignant. But Uniformity of Cell Shape > 3.5, Uniformity of Cell Size ≤ 1.5 have the opportunity to benign, etc., it is summarized in the following Table 4.3.

Decision tree model in Figure 4.9 provides an important factor used to diagnose, namely Uniformity of Cell Shape and other factors are Uniformity of Cell Size, Bland Chromatin, Single Epithelial Cell Size, Normal Nucleoli, Clump Thickness, Marginal Adhesion, and Mitoses. In the decision tree model, we are not found the Bare Nuclei. it

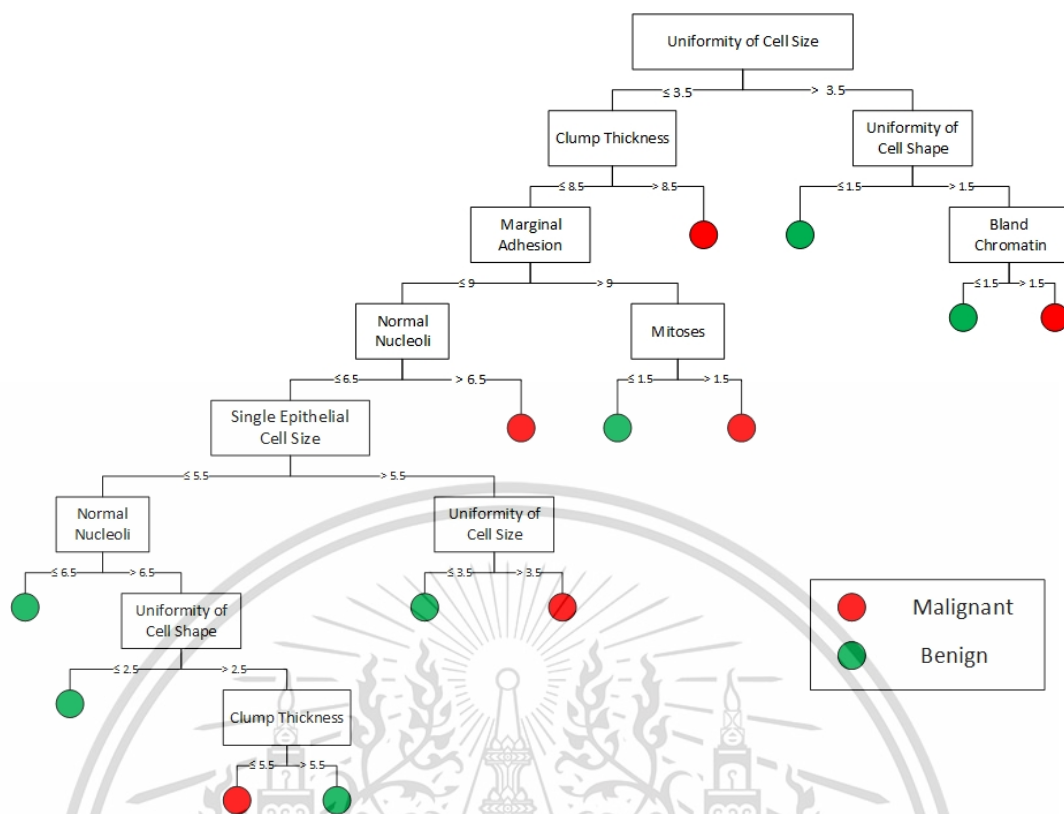


Figure 4.9: Decision tree model for testing data set by using percentage split 95:5.

means this factor impact little or none the diagnosis of breast cancer. The details of the decision tree model for the diagnosis of breast cancer are shown in Table 4.

From the confusion matrix of this model, as shown in Table 4.4, it is found that the number of the correctly classified instance, as 33 and the number of the wrongly classified instance as 2. Here TP is the number of the benign instance predicted as a benign instance. FN is the number of the benign instance predicted as a malignant instance. FP is the number of the malignant instance predicted as a benign instance. And TN is the number of the malignant instance predicted as a malignant.

It is noted that the false negative is the value of prediction error. It indicates the patients who have malignant but are predicted as benign. It is considered as the dangerous misdiagnosed which could affect the patient greatly. Also, it may cause the patient to die from a faulty diagnosis. In the data of true positive of benign, we found the value of a factor have not higher than 5. A true negative or classes of malignant of each factor are high more than 5 and less than 5 which found very few. But the false-negative appears that some factor is very high and some factors are low which may lead to a wrong prediction because each factor is not going in the same direction. Therefore, if the results of the diagnosis of the patients have the characteristic factors as the above, then the patients should be diagnosed carefully checked again. The model is developed in this paper. Therefore, taking the data for 699 instances to testing this model by the prediction from the confusion matrix shown in Table 4.4.

**Table 4.3:** The Result of Decision Tree Model for the Diagnosis Breast Cancer.

No.	UoCZ.	UoCS.	BC.	SECZ.	NN.	MA.	CT.	M.	Diagnosis.
1	$\leq 3.5$						$>8.5$		Malignant
2	$\leq 3.5$						$\leq 8.5$	$>1.5$	Malignant
3	$\leq 3.5$						$\leq 8.5$	$\leq 1.5$	Benign
4	$\leq 3.5$				$>9$	$>6.5$	$\leq 8.5$		Malignant
5	$\leq 3.5$			$>5.5$	$\leq 9$	$>6.5$	$\leq 8.5$		Malignant
6	All level			$>5.5$	$\leq 9$	$\leq 6.5$	$\leq 8.5$		Benign
7	$\leq 3.5$			$\leq 5.5$	$\leq 3.5$	$\leq 6.5$	$\leq 8.5$		Benign
8	$\leq 3.5$			$\leq 5.5$	3.5 to 9	$\leq 6.5$	$\leq 8.5$		Benign
9	$\leq 3.5$	$> 2.5$		$\leq 5.5$	3.5 to 9	$\leq 6.5$	$\leq 5.5$		Malignant
10	$\leq 3.5$	$\leq 2.5$		$\leq 5.5$	3.5 to 9	$\leq 6.5$	5.5 to 8.5		Benign
11	$> 3.5$	$\leq 1.5$	$< 1.5$				$>6.5$		Malignant
12	$> 3.5$	$\leq 1.5$	$< 1.5$				$\leq 6.5$		Benign
13	$> 3.5$	$\leq 1.5$	$< 1.5$				$\leq 6.5$		Benign

**Table 4.4:** Confusion Matrix of Decision Tree for Diagnosis of Breast Cancer by using 35 instances.

Confusion Matrix	Predicted Class		Total
	Malignant	Benign	
Pred.Malignant	21(TP)	1(FP)	22
Pred.Benign	1(FN)	12(TN)	13
Total	22	13	35

From table 4.5, the accuracy, sensitivity, and specificity are 97.0%, 96.7%, and 97.5%, respectively. We found that the rate of accuracy increases to 1.5% when using the Percentage Split method. Therefore it is reasonable to say that the method in this research, provides more accurate than that of the J48 algorithm method in [5].

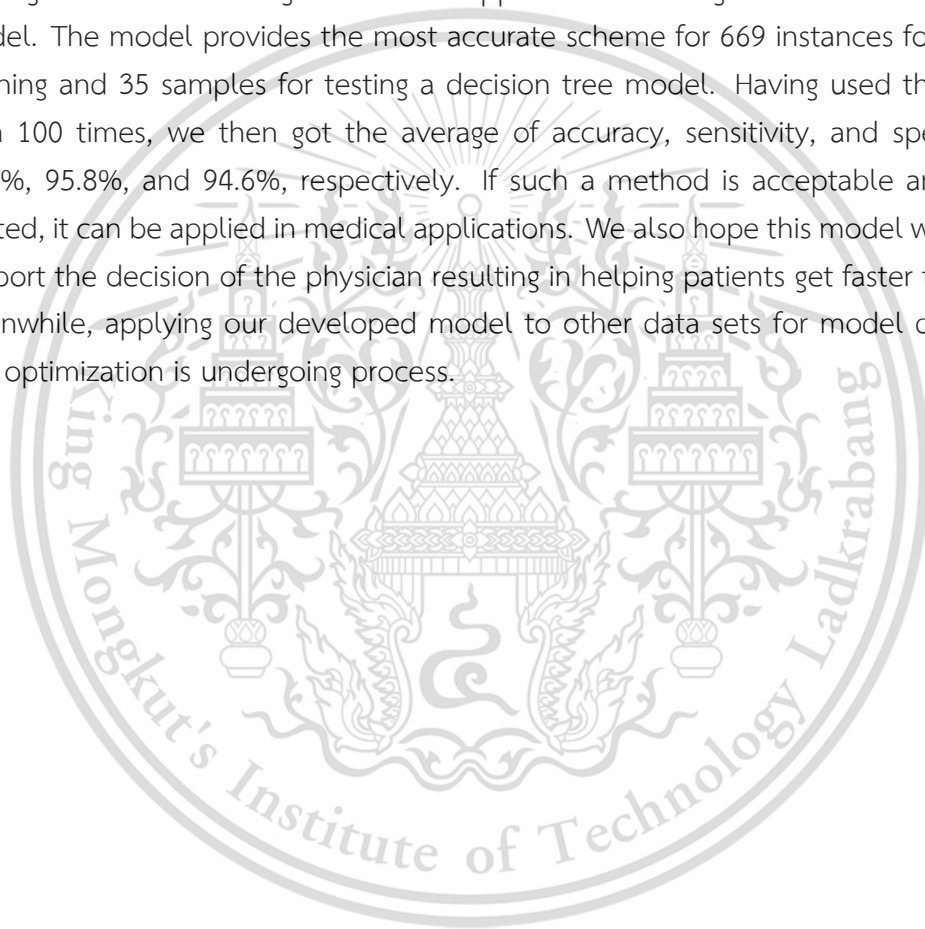
**Table 4.5:** Confusion Matrix of Decision Tree for Diagnosis of Breast Cancer by using 699 instances.

Confusion Matrix	Predicted Class		Total
	Malignant	Benign	
Pred.Malignant	443(TP)	6(FP)	447
Pred.Benign	15(FN)	235(TN)	252

## Chapter 5

### Conclusion and Discussion

This research aims to develop the model for a more accurate diagnosis based on the decision tree model and find the performance of the decision tree model that gives the best accuracy. We focused on Breast Cancer by using data from the UCI Machine Learning Data Set from Wisconsin Repository. We used a decision tree which is one of data mining technique and the Percentage Split for classifying data into the training dataset and testing dataset. We applied the C5.0 Algorithm for developing the model. The model provides the most accurate scheme for 669 instances for machine learning and 35 samples for testing a decision tree model. Having used the random data 100 times, we then got the average of accuracy, sensitivity, and specificity as 95.4%, 95.8%, and 94.6%, respectively. If such a method is acceptable and can be trusted, it can be applied in medical applications. We also hope this model would help support the decision of the physician resulting in helping patients get faster treatment. Meanwhile, applying our developed model to other data sets for model developing and optimization is undergoing process.



## References

- [1] RL. Siegel, KD. Miller, A. Jemal. Cancer Statistics, 2017. CA Cancer J Clin., 7-30, 2017.
- [2] J. Han, M. Kamber, J. Pei, Data mining, Concepts and Techniques 3rd ed., USA, 2011.
- [3] M.Shouman, T. Turner, R. Stocker, Using Decision Tree for Diagnosing Heart Disease Patients. Proceedings of the 9th Australasian Data Mining Conference, Ballarat, Australia, 2011.
- [4] M. Bramer, Principles of Data Mining (Third Edition), Undergraduate Topics in Computer Science, Springer Verlag London Ltd., 2016.
- [5] R. Sumbaly, N. Vishusri, S. Jeyalatha, Diagnosis of Breast Cancer using Decision Tree Data Mining Technique. International Journal of Computer Applications Volume 98-No.10, 16-24, 2014.
- [6] D. Muntham, L. Ingsrisawang, An Application of Decision Tree Algorithms for Diagnosis of the Respirator System: A Case Study of Pranakorn Sri Ayudthaya Hospital. Journal of Health Systems Research Vol.4 No.1, 72-81, 2010.
- [7] S. Auwatanamongkol. Data Mining. Bangkok Block Ltd., Part, Bangkok, 31-41, 2014.
- [8] S. Dharm, Naveen C.and Jully S., Analysis of Data Mining Classification with Decision tree Technique, Global Journal of Computer Science and Technology Software & Data Engineering, Vol. 13 Issue 13, 2013.
- [9] A. N. Richter, & Khoshgoftaar, T. M., A review of statistical and machine learning methods for modeling cancer risk using structured clinical data. Artificial Intelligence in Medicine., 2018.
- [10] เอื้อวัฒนามงคล สุรพงษ์, การทำเหมืองข้อมูล, สำนักพิมพ์สถาบันจิตพัฒนาบริหารศาสตร์, 2010.
- [11] สีนสมบูรณ์ทอง สายชล, การทำเหมืองข้อมูล, จามจุรีโปรดักท์, 2009.
- [12] Sokolova M., Lapalme G., Performance Measures in Classification of Human Communications.,Advances in Artificial Intelligence. AI 2007. Lecture Notes in Computer Science, vol 4509. Springer,Berlin, Heidelberg, 2007.
- [13] George D. M., Andriana P., Machine Learning in Medical Applications, Springer-Verlag Berlin Heidelberg, 300-307, 2001.
- [14] Frank E., Hall M., Holmes G., Kirkby R., Pfahringer B., Witten H. I., Trigg L., Weka-A Machine Learning Workbench for Data Mining, In: Maimon O., Rokach L. (eds) Data Mining and KnowledgeDiscovery Handbook. Springer, Boston, MA, 2010.

- [15] Frank, E., Hall, M., Trigg, L., Holmes, G., Witten, I. H., Data mining in bioinformatics using Weka. *Bioinformatics*, 20(15), 2479–2481., 2004.
- [16] Therneau T., Atkinson B., An Introduction to Recursive Partitioning Using the RPART Routines., <https://cran.r-project.org/web/packages/rpart/vignettes/longintro.pdf>., 2018.
- [17] Meyer D., Support Vector Machines—the Interface to libsvm in package e1071, <https://cran.r-project.org/web/packages/e1071/vignettes/svmdoc.pdf>., 2018.
- [18] Asri, H., Mousannif, H., Moatassime, H. A., Noel, T., Using Machine Learning Algorithms for Breast Cancer Risk Prediction and Diagnosis. *Procedia Computer Science*, 83, 1064–1069., 2016.
- [19] Rau, H.-H., Hsu, C.-Y., Lin, Y.-A., Atique, S., Fuad, A., Wei, L.-M., Hsu, M.-H., Development of a web-based liver cancer prediction model for type II diabetes patients by using an artificial neural network. *Computer Methods and Programs in Biomedicine*, 125, 58–65., 2016.
- [20] Cirkovic, B. R. A., Cvetkovic, A. M., Ninkovic, S. M., Filipovic, N. D., Prediction models for estimation of survival rate and relapse for breast cancer patients., 2015 IEEE 15th International Conference on Bioinformatics and Bioengineering (BIBE)., 2015.
- [21] Ahmad LG, Eshlaghy AT, Poorebrahimi A, Ebrahimi M, Razavi AR., Using Three Machine Learning Techniques for Predicting Breast Cancer Recurrence. *J Health Med Inform*, 2013.
- [22] Montazeri, M., Montazeri, M., Montazeri, M., Beigzadeh, A., Machine learning models in breast cancer survival prediction. *Technology and Health Care*, 24(1), 2016.
- [23] S. Perveen., M. Shahbaz, K. Keshajee, A. Guergachi, A Systematic Machine Learning Based Approach for the Diagnosis of Non-Alcoholic Fatty Liver Disease Risk and Progression., *Scientific reports*, 2018.
- [24] B. R. A. Cirkovic, Cvetkovic, A. M., Ninkovic, S. M., & Filipovic, N. D., Prediction models for estimation of survival rate and relapse for breast cancer patients., 2015 IEEE 15th International Conference on Bioinformatics and Bioengineering (BIBE)., 2015.
- [25] Hamsagayathri, P., & Sampath, P. Priority based decision tree classifier for breast cancer detection. 2017 4th International Conference on Advanced Computing and Communication Systems (ICACCS)., 2017.
- [26] Chaurasia, V., Pal, S., & Tiwari, B. Prediction of benign and malignant breast cancer using data mining techniques. *Journal of Algorithms & Computational Technology*, 119–126., 2018.

This material is intended for educational use only, not allowed for commercial use.

Forbidden to modify the content, and cite the document when use.

[27] Dr. William H. Wolberg, Olvi Mangasarian, University of Wisconsin Hospitals Madison, Wisconsin, USA, [https://archive.ics.uci.edu/ml/datasets/breast+cancer+wisconsin+\(original\)](https://archive.ics.uci.edu/ml/datasets/breast+cancer+wisconsin+(original))



This material is reserved for educational use only, not allowed for commercial use.  
Forbidden to modify the content, and cite the document when use.



This material is reserved for educational use only, not allowed for commercial use.  
Forbidden to modify the content, and cite the document when use.



No	Color	Uniformity of color	Uniformity of texture	Internal Adhesion	Impact Resistance	Flame Retardant	Field Chromatic	Resistance to aging	Class
81	1137156	3	3	3	3	3	3	3	1. Plastic
82	1137358	4	4	4	4	4	4	4	1. Plastic
83	1143978	5	2	1	1	2	1	3	1. Plastic
84	1147044	2	3	3	3	3	3	3	1. Plastic
85	1147698	3	5	2	8	8	9	7	10. 7. Metallurgical
86	1147400	2	10	2	1	10	2	2	10. 7. Metallurgical
87	1148778	3	3	3	3	3	3	3	1. Plastic
88	1148673	3	0	0	0	0	10	3	1. Plastic
89	1150554	4	3	3	3	3	3	3	1. Plastic
90	1150546	2	2	2	2	2	2	2	1. Plastic
91	1150252	2	2	2	2	2	2	2	1. Plastic
92	1150340	0	0	0	0	0	0	0	1. Plastic
93	1157734	4	5	2	2	2	2	2	1. Plastic
94	1158242	1	1	1	1	1	1	1	1. Plastic
95	1160436	3	3	3	3	3	3	3	1. Plastic
96	1160804	3	3	3	3	3	3	3	1. Plastic
97	1160917	2	1	1	2	2	1	1	1. Plastic
98	1161374	3	3	3	3	3	3	3	1. Plastic
99	1160954	2	0	0	0	0	0	0	1. Plastic
100	1160230	2	2	2	2	2	2	2	1. Plastic
101	1160250	10	0	0	0	0	0	0	1. Plastic
102	1167429	3	2	4	4	3	5	3	5. 1. Metallurgical
103	1167475	4	2	4	4	2	5	3	5. 1. Metallurgical
104	1168358	8	2	3	3	6	3	2	2. Metallurgical
105	1168736	10	10	10	10	10	1	8	8. Metallurgical
106	1169045	7	3	4	4	3	3	3	7. Metallurgical
107	1170010	10	10	10	10	10	2	2	10. Metallurgical
108	1170676	1	8	8	10	8	10	3	1. Metallurgical
109	1171700	1	1	1	1	2	1	3	1. Metallurgical
110	1171710	0	0	0	0	0	0	0	1. Metallurgical
111	1171759	1	5	1	2	2	3	3	1. Metallurgical
112	1171860	0	0	0	0	0	0	0	1. Metallurgical
113	1172152	10	3	3	10	2	10	3	1. Metallurgical
114	1173216	10	10	10	10	10	10	10	1. Metallurgical
115	1173325	2	2	2	2	2	2	2	1. Metallurgical
116	1173282	3	3	3	3	3	5	3	1. Metallurgical
117	1173280	3	3	3	3	3	3	3	1. Metallurgical
118	1173493	3	3	3	3	3	3	3	1. Metallurgical
119	11734168	1	1	1	1	1	1	1	1. Metallurgical
120	11733664	5	4	4	4	4	4	4	1. Metallurgical
121	11734857	4	3	3	3	3	3	3	1. Metallurgical
122	1173637	4	2	2	2	2	2	2	1. Metallurgical
123	1174151	10	10	10	10	10	10	10	1. Metallurgical
124	1174428	5	3	6	5	5	5	5	1. Metallurgical
125	1175932	5	4	6	7	8	7	10	1. Metallurgical
126	1176406	1	1	1	1	1	1	1	1. Metallurgical
127	1176881	7	5	3	7	4	10	2	5. Metallurgical
128	1176044	4	4	4	4	4	4	4	1. Metallurgical
129	1173398	8	1	3	8	1	10	3	1. Metallurgical
130	1177512	1	1	1	1	1	1	1	1. Metallurgical
131	1178380	3	1	3	2	1	2	3	1. Metallurgical
132	1179050	2	2	2	2	2	2	2	1. Metallurgical
133	1180054	0	10	0	10	0	10	0	1. Metallurgical
134	1180223	0	0	0	0	0	0	0	1. Metallurgical
135	1180038	3	3	3	3	3	3	3	1. Metallurgical
136	1181256	2	2	2	2	2	2	2	1. Metallurgical
137	1182404	4	3	3	3	3	3	3	1. Metallurgical
138	1182600	4	3	3	3	3	3	3	1. Metallurgical
139	1183200	4	3	3	3	3	3	3	1. Metallurgical
140	1183360	1	1	1	1	1	1	1	1. Metallurgical
141	1183380	3	3	3	3	3	3	3	1. Metallurgical
142	1183911	2	1	1	2	1	1	1	1. Metallurgical
143	1183963	3	3	3	3	3	3	3	1. Metallurgical
144	1184182	0	0	0	0	0	0	0	1. Metallurgical
145	1184242	2	2	2	2	2	2	2	1. Metallurgical
146	1184865	3	3	3	3	3	3	3	1. Metallurgical
147	1185603	3	4	5	2	6	4	1	1. Metallurgical
148	1185610	1	1	1	1	3	2	2	1. Metallurgical
149	1187857	3	1	3	1	8	1	5	2. Metallurgical
150	1187805	8	8	8	8	10	8	3	1. Metallurgical
151	1188472	4	4	4	4	4	4	4	1. Metallurgical
152	1189340	7	2	2	2	2	2	4	1. Metallurgical
153	1189260	10	10	0	0	0	0	0	1. Metallurgical
154	1190334	4	2	2	2	2	2	2	1. Metallurgical
155	1190402	0	0	0	0	0	0	0	1. Metallurgical
156	1192336	1	1	1	1	1	1	1	1. Metallurgical
157	1193083	1	2	2	2	2	2	2	1. Metallurgical
158	1193210	2	1	1	1	2	1	2	1. Metallurgical
159	1193400	1	1	1	1	1	1	1	1. Metallurgical
160	1196743	4	4	10	4	10	4	10	1. Metallurgical
161	1196443	10	1	1	1	10	3	2	1. Metallurgical
162	1197000	4	1	4	4	4	4	4	1. Metallurgical
163	1202270	0	0	0	0	0	0	0	1. Metallurgical
164	1197980	2	2	2	2	2	2	2	1. Metallurgical
165	1197510	5	1	1	1	2	3	1	1. Metallurgical
166	1192338	4	5	5	5	5	5	5	1. Metallurgical
167	1192992	5	6	6	6	6	6	6	1. Metallurgical
168	1198178	10	8	10	10	6	1	3	10. Metallurgical
169	1198641	3	1	1	1	2	1	3	1. Metallurgical
170	1199100	7	7	7	7	7	7	7	1. Metallurgical
171	1199231	3	1	1	1	2	1	1	1. Metallurgical
172	1199960	1	1	1	1	2	1	1	1. Metallurgical
173	1200772	1	1	1	1	2	1	1	1. Metallurgical
174	1200847	0	10	10	10	10	10	10	1. Metallurgical
175	1206632	0	0	0	0	0	0	0	1. Metallurgical
176	1206952	0	0	0	0	0	0	0	1. Metallurgical
177	1201834	0	0	0	0	0	0	0	1. Metallurgical
178	1201925	0	10	10	10	10	10	10	1. Metallurgical
179	1202135	4	3	3	3	3	3	3	1. Metallurgical
180	1202800	4	4	4	4	4	4	4	1. Metallurgical

This material is reserved for educational use only, not allowed for commercial use.

Forbidden to modify the content, and cite the document when use.









No	Completion	Uniformity of	Adhesion	Initial	Final	Chromatic	Mass	
591	1294552	4	0	0	0	0	Benign	
592	1295208	4	5	0	0	0	Benign	
593	1296000	10	7	0	0	0	Benign	
594	1311825	5	1	0	0	0	Benign	
595	1315006	4	8	0	0	10	Benign	
596	1320000	0	0	0	0	0	Benign	
597	1321000	0	0	0	0	0	Benign	
598	1333000	0	0	0	0	0	Benign	
599	1333000	0	0	0	0	0	Benign	
600	1334000	0	0	0	0	0	Benign	
601	1328700	0	0	0	0	0	Benign	
602	1244422	0	0	0	0	0	Benign	
603	1350568	4	0	0	0	0	Benign	
604	1352663	5	4	0	0	0	Benign	
605	188326	5	2	0	0	10	Benign	
606	852000	10	0	0	0	0	Benign	
607	350000	4	1	0	0	0	Benign	
608	411433	0	0	0	0	0	Benign	
609	209268	0	0	0	0	0	Benign	
610	202792	0	0	0	0	0	Benign	
611	352350	0	0	0	0	0	Benign	
612	879200	0	0	0	0	0	Benign	
613	858025	0	0	0	0	0	Benign	
614	1035535	0	0	0	0	0	Benign	
615	1031699	0	0	0	0	0	Benign	
616	1041042	0	0	0	0	0	Benign	
617	1042752	0	0	0	0	0	Benign	
618	1037000	0	0	0	0	0	Benign	
619	1061900	0	0	0	0	0	Benign	
620	1072800	0	0	0	0	0	Benign	
621	1083817	0	0	0	0	0	Benign	
622	1056332	0	0	0	0	0	Benign	
623	1140000	0	0	0	0	0	Benign	
624	1169550	0	0	0	0	0	Benign	
625	1130000	0	0	0	0	0	Benign	
626	1182500	0	0	0	0	0	Benign	
627	1101000	0	0	0	0	0	Benign	
628	1190000	0	0	0	0	0	Benign	
629	1215000	0	0	0	0	0	Benign	
630	1233000	0	0	0	0	0	Benign	
631	1235000	0	0	0	0	0	Benign	
632	1238000	0	0	0	0	0	Benign	
633	1238000	0	0	0	0	0	Benign	
634	1253000	0	0	0	0	0	Benign	
635	1254000	0	0	0	0	0	Benign	
636	1260000	0	0	0	0	0	Benign	
637	1268900	10	10	0	0	10	10	Benign
638	1270000	0	0	0	0	0	0	Benign
639	1272000	0	0	0	0	0	0	Benign
640	1277000	0	0	0	0	0	0	Benign
641	1280000	0	0	0	0	0	0	Benign
642	1266000	0	0	0	0	0	0	Benign
643	1283000	0	0	0	0	0	0	Benign
644	1284000	0	0	0	0	0	0	Benign
645	1299000	0	0	0	0	0	0	Benign
646	1303000	0	0	0	0	0	0	Benign
647	1311000	0	0	0	0	0	0	Benign
648	1300000	0	0	0	0	0	0	Benign
649	1315000	0	0	0	0	0	0	Benign
650	1310000	0	0	0	0	0	0	Benign
651	1305000	0	0	0	0	0	0	Benign
652	1305000	0	0	0	0	0	0	Benign
653	1328000	0	0	0	0	0	0	Benign
654	1324000	0	0	0	0	0	0	Benign
655	1325000	0	0	0	0	0	0	Benign
656	1330000	0	0	0	0	0	0	Benign
657	1330000	0	0	0	0	0	0	Benign
658	1338000	0	0	0	0	0	0	Benign
659	1340000	0	0	0	0	0	0	Benign
660	1340000	0	0	0	0	0	0	Benign
661	1337000	0	0	0	0	0	0	Benign
662	1337000	0	0	0	0	0	0	Benign
663	1345000	0	0	0	0	0	0	Benign
664	1343000	0	0	0	0	0	0	Benign
665	1345000	0	0	0	0	0	0	Benign
666	1343000	0	0	0	0	0	0	Benign
667	1347000	0	0	0	0	0	0	Benign
668	1348000	0	0	0	0	0	0	Benign
669	1346000	0	0	0	0	0	0	Benign
670	1350000	0	0	0	0	0	0	Benign
671	1350000	0	0	0	0	0	0	Benign
672	1335000	0	0	0	0	0	0	Benign
673	1330000	0	0	0	0	0	0	Benign
674	1328000	0	0	0	0	0	0	Benign
675	1325000	0	0	0	0	0	0	Benign
676	1325000	0	0	0	0	0	0	Benign
677	1305000	0	0	0	0	0	0	Benign
678	1368000	0	0	0	0	0	0	Benign
679	1368000	0	0	0	0	0	0	Benign
680	1368000	0	0	0	0	0	0	Benign
681	1370000	0	0	0	0	0	0	Benign
682	1371000	0	0	0	0	0	0	Benign
683	1373000	0	0	0	0	0	0	Benign
684	1369000	0	0	0	0	0	0	Benign
685	1369000	0	0	0	0	0	0	Benign
686	1369000	0	0	0	0	0	0	Benign
687	1369000	0	0	0	0	0	0	Benign
688	1362000	0	0	0	0	0	0	Benign
689	1362000	0	0	0	0	0	0	Benign
690	1362000	0	0	0	0	0	0	Benign
691	1362000	0	0	0	0	0	0	Benign
692	1362000	0	0	0	0	0	0	Benign
693	1362000	0	0	0	0	0	0	Benign
694	1362000	0	0	0	0	0	0	Benign
695	1362000	0	0	0	0	0	0	Benign
696	1362000	0	0	0	0	0	0	Benign
697	1362000	0	0	0	0	0	0	Benign
698	1362000	0	0	0	0	0	0	Benign
699	1362000	0	0	0	0	0	0	Benign

This material is reserved for educational use only, not allowed for commercial use.

Forbidden to modify the content, and cite the document when use.

## Author Biography

Name	Miss Jiraphorn Charoenpong
Date of Birth	17 September 1994
Address	99/212, Kanjanapisek Road, Prawet, Prawet District Bangkok, 10250
Education	2017 Bachelor of Science in Applied Mathematics GPA 2.86 King Mongkut's Institute of Technology Ladkrabang 2019 Master of Science in Applied Mathematics GPA 3.32 King Mongkut's Institute of Technology Ladkrabang
Scholarship(s)	Assistant Researcher by Centre of Excellence in Mathematics Assistant Researcher by Thailand Center of Excellence in Physics
Academic Publication(s)	1. Jiraporn C., Busayamas P., Somkid A., Wannapong T. and Narin N. "A Comparison of Machine Learning Algorithms and their Applications", International Journal of Simulation Systems, Science and Technology, Volume 20, Number 4, August 2019. DOI 10.5013/IJSSST.a.20.04.08