

**FEATURE SELECTION USING FAST CORRELATION-BASED FILTER  
(FCBF) FOR LEAF CLASSIFICATION**

LUA XIN LIN

This report is submitted in partial fulfilment of the requirements for the Bachelor of  
Computer Science (Artificial Intelligence)

FACULTY OF INFORMATION AND COMMUNICATION TECHNOLOGY  
UNIVERSITI TEKNIKAL MALAYSIA MELAKA  
2015

## BORANG PENGESAHAN STATUS TESIS\*

JUDUL:               **FEATURE SELECTION USING FAST CORRELATION-BASED FILTER (FCBF) FOR LEAF CLASSIFICATION**

SESI

PENGAJIAN:   **2014/2015**

Saya **LUA XIN LIN** mengaku membenarkan tesis (PSM/Sarjana/Doktor Falsafah) ini disimpan di Perpustakaan Fakulti Teknologi Maklumat dan Komunikasi dengan syarat-syarat kegunaan seperti berikut:

1. Tesis dan projek adalah hakmilik Universiti Teknikal Malaysia Melaka.
2. Perpustakaan Fakulti Teknologi Maklumat dan Komunikasi dibenarkan membuat
3. salinan untuk tujuan pengajian sahaja.
4. Perpustakaan Fakulti Teknologi Maklumat dan Komunikasi dibenarkan membuat
5. salinan tesis ini sebagai bahan pertukaran antara institusi pengajian tinggi.
6. \*\* Sila tandakan (/)

|       |              |  |
|-------|--------------|--|
| _____ | SULIT        | (Mengandungi maklumat yang berdarjah keselamatan atau kepentingan Malaysia seperti yang termaktub di dalam AKTA RAHSIA RASMI 1972) |
| _____ | TERHAD       | (Mengandungi maklumat TERHAD yang telah ditentukan oleh organisasi/badan di mana penyelidikan dijalankan)                          |
| _/_   | TIDAK TERHAD |  |

---

(TANDATANGAN PENULIS)

---

(TANDATANGAN PENYELIA)

**ALAMAT TETAP:** 49 Jln BU4/9  
Bandar Utama 47800 Petaling Jaya,  
Selangor Darul Ehsan, Malaysia

---

(NAMA PENYELIA)

**Tarikh:**

**Tarikh:**

CATATAN:

\* Tesis dimaksudkan sebagai Laporan Akhir Projek Sarjana Muda (PSM)

\*\* Jika tesis ini SULIT atau TERHAD, sila lampirkan surat daripada pihak berkuasa.

## DECLARATION

I hereby declare that this project report entitled

**FEATURE SELECTION USING FAST CORRELATION-BASED FILTER  
(FCBF) FOR LEAF CLASSIFICATION**

is written by me and is my own effort and that no part has been plagiarized without  
citation.

STUDENT: \_\_\_\_\_ DATE: \_\_\_\_\_

(LUA XIN LIN)

SUPERVISOR: \_\_\_\_\_ DATE: \_\_\_\_\_

(DR AZAH KAMILAH BINTI DRAMAN @ MUDA)

## DEDICATION

This work is dedicated to my family and friends. To my parents, for your unwavering support—mom and dad, for always nagging at me to get on with my work. My brother, Alex, for your sarcasm and wit—really, what are siblings for if not some friendly squabbling? To my dear friend Zoe, for not begrudging my late night calls, patiently listening to all my woes and the many words of encouragement—you're my rock. To my supervisor, Dr Azah Kamilah bt. Draman@Muda, for making all this possible.

## ACKNOWLEDGEMENT

I'd like to thank my supervisor, Dr Azah Kamilah bt Draman@Muda, who has been more than generous with her expertise and time. Your guidance is priceless.

I'd also like to thank Dr Choo Yun Huoy, for willingly imparting advice and information. I immensely enjoyed your Machine Learning lectures; those lessons have aided me well in the process of this work.

## ABSTRACT

Imagine taking a forest trail hike, where you see many interesting plants. How would you differentiate one plant from another based on its leaves? Do you base your decision on its color? Or its texture? Or size? The thing is, one single leaf could give you a lot of information, from color to texture to size, so on and so forth. But the problem is: which features to focus on so as to make a good decision?

The identification of a feature subset that best represents a class so as to build a strong predictive model is still an issue that researchers are still working on solving. This paper focuses on feature selection, where the Fast Correlation-Based Filter is compared to the Correlation-based Feature Selector.

Four datasets were retrieved from the UCI Machine Learning Repository. The SVM classifier was used to build a predictive model based on the feature subset. The experiments were done using the weka machine learning tool, with the evaluation criterion being predictive accuracy, kappa measure, and time taken to build a model.

## ABSTRAK

Bayangkan, anda sedang trek di hutan, di mana anda melihat banyak tumbuhan yang menarik. Bagaimana anda membezakan satu tumbuhan dari yang lain berdasarkan daun? Adakah anda mendasarkan keputusan anda pada warnanya? Atau tekstur? Atau saiz? Masalahnya, satu daun tunggal boleh memberikan banyak maklumat, dari segi warna, tekstur, saiz dan sebagainya. Tetapi masalahnya ialah: ciri-ciri manakah yang harus diberi tumpuan supaya keputusan yang baik dapat dilakukan?

Pengenalpastian subset 'feature', ataupun ciri yang terbaik mewakili kelas bagi membina model ramalan yang kuat masih merupakan isu yang penyelidik masih berusaha untuk menyelesaikan. Tesis ini memberi tumpuan kepada pemilihan 'feature', di mana Fast Correlation Based Filter dibandingkan dengan ciri Correlation-Based Filter.

Empat set data telah diperolehi dari UCI Machine Learning Repository. Classifier SVM telah digunakan untuk membina model ramalan berdasarkan subset 'feature' tersebut. Eksperimen telah dilakukan dengan menggunakan alat pembelajaran mesin weka, dengan kriteria penilaian yang berikut: ketepatan ramalan, statistik kappa, dan masa yang diambil untuk membina model.

## TABLE OF CONTENTS

| CHAPTER    | SUBJECT                       | PAGE |
|------------|-------------------------------|------|
|            | DECLARATION OF THESIS STATUS  | i    |
|            | DECLARATION                   | iii  |
|            | DEDICATION                    | iv   |
|            | ACKNOWLEDGEMENTS              | v    |
|            | ABSTRACT                      | vi   |
|            | ABSTRAK                       | vii  |
|            | TABLE OF CONTENTS             | viii |
|            | LIST OF TABLES                | xii  |
|            | LIST OF FIGURES               | xiv  |
|            | LIST OF ABBREVIATIONS         | xvi  |
| CHAPTER I  | INTRODUCTION                  |      |
|            | 1.1 Introduction              | 1    |
|            | 1.2 Problem Statement         | 2    |
|            | 1.3 Objectives                | 3    |
|            | 1.4 Scope                     | 3    |
|            | 1.5 Project Significance      | 4    |
|            | 1.6 Expected Result           | 4    |
|            | 1.7 Conclusion                | 5    |
| CHAPTER II | LITERATURE REVIEW             | 6    |
|            | 2.1 Introduction              | 6    |
|            | 2.2 Plant Leaf Identification | 7    |



|           |   |    |
|-----------|---|----|
| 2.2.1     | Feature Extraction  | 8  |
| 2.2.2     | Feature Selection   | 9  |
| 2.2.2.1   | Filter Model  | 12 |
| 2.2.2.2   | Wrapper Model   | 12 |
| 2.2.3     | Feature Relevance   | 13 |
| 2.2.4     | Feature Redundancy  | 14 |
| 2.2.5     | Existing Feature Selection<br>Techniques  | 14 |
| 2.2.5.1   | In Plant Classification   | 15 |
| 2.2.5.1.1 | Principle Component<br>Analysis   | 15 |
| 2.2.5.1.2 | Genetic Algorithm   | 15 |
| 2.2.5.2   | Other Techniques  | 16 |
| 2.2.5.2.1 | Particle Swarm Optimization   | 16 |
| 2.3       | Existing Work on Plant Leaf<br>Identification   | 17 |
| 2.3.1     | A leaf Recognition Algorithm for<br>Plant Classification using PNN                        | 17 |
| 2.3.2     | A Rapid Flower/Leaf<br>Recognition System   | 19 |
| 2.3.3     | Experiments of Zernike moment<br>for Leaf Identification                                  | 20 |
| 2.3.4     | An Optimal Feature Subset<br>Selection using Genetic<br>Algorithm for Leaf Identification | 22 |
| 2.3.5     | Texture Feature and k-Nearest<br>Neighbor in Classification of<br>Flower Images           | 23 |
| 2.4       | Discussion  | 24 |
| 2.4.1     | Support Vector Classifier   | 25 |

|             |  |    |
|-------------|--|----|
| 2.4.2       | Evaluation Criteria                    | 29 |
| 2.4.2.1     | Correlation                            | 30 |
| 2.4.2.2     | Symmetrical Uncertainty<br>Coefficient | 31 |
| 2.4.3       | Correlation-based Feature<br>Selector  | 32 |
| 2.4.3.1     | Best-First Search                      | 33 |
| 2.4.3.2     | Particle Swarm Optimization            | 33 |
| 2.4.4       | Fast Correlation Based Filter          | 34 |
| 2.5         | Requirements                           | 36 |
| 2.5.1       | Software                               | 36 |
| 2.5.1.1     | Weka                                   | 37 |
| 2.5.1.2     | Microsoft Office 2010                  | 37 |
| 2.5.1.3     | MATLAB                                 | 37 |
| 2.5.2       | Hardware                               | 37 |
| 2.6         | Project Schedule and Milestone         | 37 |
| 2.7         | Conclusion                             | 38 |
| CHAPTER III | METHODOLOGY                            | 39 |
| 3.1         | Introduction                           | 39 |
| 3.2         | Data Representation                    | 40 |
| 3.2.1       | Data set 1                             | 41 |
| 3.2.2       | Dataset 2                              | 41 |
| 3.3         | Discretization                         | 42 |
| 3.3.1       | Unsupervised discretization            | 44 |
| 3.3.2       | Supervised discretization              | 44 |
| 3.4         | Data Sampling                          | 46 |
| 3.4.1       | Percentage Split                       | 46 |
| 3.4.2       | Use Test Set                           | 47 |
| 3.4.3       | k-fold Cross-validation                | 48 |

|  |           |
|--|-----------|
| 3.5 Sequential Minimal Optimization (SMO)                              | 49        |
| 3.5.1 Solving the Quantization Problem (QP) through Analytical Methods | 49        |
| 3.5.2 Heuristics to choose the Correct Multiplier to Optimize          | 51        |
| 3.5.3 Computing threshold value b                                      | 52        |
| 3.6 Performance Evaluation   | 52        |
| 3.6.1 Classifier Accuracy  | 52        |
| 3.6.2 Time taken to generate a model                                   | 53        |
| 3.6.3 Kappa statistics   | 55        |
| 3.7 Experiment Approach  | 55        |
| 3.7.1 Weka   | 55        |
| 3.7.2 MATLAB   | 56        |
| 3.7.3 The approach   | 56        |
| 3.8 Conclusion   |           |
| <b>CHAPTER IV RESULTS AND ANALYSIS</b>                                 | <b>57</b> |
| 4.1 Results  | 57        |
| 4.1.1 Discretized vs. Undiscretized                                    | 57        |
| 4.1.1.1 Dataset1   | 57        |
| 4.1.1.2 Dataset2   | 58        |
| 4.1.2 Discretized +FS vs Undiscretized+FS                              | 61        |
| 4.1.2.1 Dataset 1  | 61        |
| 4.1.2.1.1. Margin  | 61        |
| 4.1.2.1.2. Shape   | 63        |
| 4.1.2.1.3. Texture   | 64        |
| 4.1.2.2. Dataset 2   | 65        |

|  |    |
|--|----|
| 4.1.3. Supervised Discretization + FS vs.<br>Unsupervised Discretization | 71 |
| 4.1.3.1. Margin  | 71 |
| 4.1.3.2. Shape   | 72 |
| 4.1.3.3. Texture   | 72 |
| 4.1.3.4. dataset 2   | 72 |
| 4.2 Analysis & Discussion  | 73 |
| 4.2.1 Introduction   | 73 |
| 4.2.2 Discussion   | 74 |
| 4.3 Conclusion.  | 78 |
| <br>CHAPTER V  |    |
| CONCLUSION   | 79 |
| 5.1 Observation on Weaknesses and<br>Strength                            | 79 |
| 5.2 Proposition for Improvement  | 79 |
| 5.3 Project Contribution   | 80 |
| 5.4 Conclusion   | 80 |
| <br>REFERENCES   | 82 |
| BIBLIOGRAPHY   | 86 |
| APPENDIX   | 87 |

## LIST OF TABLES

| TABLE | TITLE  | PAGE |
|-------|--|------|
| 4.1   | Undiscretized Dataset  | 58   |
| 4.2   | Supervised discretization  | 59   |
| 4.3   | Unsupervised Discretization  | 59   |
| 4.4   | Undiscretized, supervised discretized,<br>unsupervised-discretization  | 59   |
| 4.5   | FS on Undscretized vs Undiscretized<br>margin dataset  | 62   |
| 4.6   | FS on Undscretized vs Undiscretized<br>shape dataset   | 64   |
| 4.7   | FS on Undscretized vs Undiscretized<br>texture   | 65   |
| 4.8   | FS on Undscretized vs Undiscretized<br>dataset2  | 66   |
| 4.9   | comparison of SVM on full margin<br>dataset preprocessed with supervised<br>discretization and feature subset<br>generated from dataset preprocessed with<br>supervised discretization | 72   |
| 4.10  | comparison of SVM on full shape dataset<br>preprocessed with supervised<br>discretization and feature subset<br>generated from dataset preprocessed with<br>supervised discretization  | 73   |
| 4.11  | comparison of SVM on full dataset<br>preprocessed with supervised  | 73   |

|      |  |    |
|------|--|----|
|      | discretization and feature subset<br>generated from dataset preprocessed with<br>supervised discretization   |    |
| 4.12 | comparison of SVM on full dataset2<br>preprocessed with supervised<br>discretization and feature subset<br>generated from dataset preprocessed with<br>supervised discretization | 74 |

## LIST OF FIGURES

| DIAGRAM | TITLE   | PAGE |
|---------|---|------|
| 2.1     | Process of Plant Leaf Classification  | 7    |
| 2.2     | 4 basic steps in Feature Selection (Dash & Liu, 2005)   | 10   |
| 2.3     | process flow of the filter model  | 12   |
| 2.4     | Wrapper method Process flow   | 13   |
| 2.5     | A Leaf Recognition Algorithm for Plant Classification Using Probabilistic Neural Network (Wu et. al, 2007)            | 17   |
| 2.6     | Artificial Neural Network   | 18   |
| 2.7     | framework of flower/leaf recognition system (Qi, et. al, 2012)  | 19   |
| 2.8     | System using PNN (Kadir, et.al, 2003)   | 21   |
| 2.9     | system using Distance Measure (Kadir, et.al, 2003)  | 22   |
| 2.10    | framework of a plant classification system using SVM using GA-based subset selection (Valliammal and Subbaraya, 2014) | 23   |
| 2.11    | Optimal Separating Hyper plane (Gunn, 1998)   | 26   |
| 2.12    | Canonical Hyperplane (Gunn,1998)  | 27   |
| 2.13    | Constraining the Canonical Hyperspace (Gunn, 1998)  | 28   |

|      |  |    |
|------|--|----|
| 2.14 | Standard PSO algorithm                                     | 34 |
| 2.15 | pseudocode o the FCBF algorithm                            | 36 |
| 3.1  | arrangement of data  | 40 |
| 3.2  | Percentage Split   | 47 |
| 3.3  | use test set   | 47 |
| 3.4  | Cross-validation   | 48 |
| 3.5  | Output model   | 48 |
| 3.6  | Langrangian constraints (Platt, 1999)                      | 49 |
| 3.7  | confusion matrix   | 53 |
| 3.8  | Kappa Statistics Interpretation (Landis<br>and Koch, 1977) | 55 |



**LIST OF ABBREVIATIONS**

|      |    |                                     |
|------|----|-------------------------------------|
| FS   | -- | Feature Selection                   |
| PSO  | -- | Particle Swarm Optimization         |
| BFS  | -- | Best First Search                   |
| SVM  | -- | Support Vector Machine              |
| GCLM | -- | Grey-Level Co-concurrence Matrix    |
| SMO  | -- | Sequential Minimal Optimization     |
| KKT  | -- | Karush-Kuhn Tucker                  |
| FCBF | -- | Fast Correlation-Based Filter       |
| CFS  | -- | Correlation-Based Feature Selection |
| PNN  | -- | Probabilistic Neural Networks       |
| kNN  | -- | k-Nearest Neighbor                  |
| GA   | -- | Genetic Algorithm                   |
| PCA  | -- | Principle Component Analysys        |

## **CHAPTER 1**

### **INTRODUCTION**

#### **1.1 Introduction**

One of the upcoming fields in need of research and development is the automatic recognition of plant species—plant classification. This is due to the vital role plants play in the natural circle of life: the conversion of carbon dioxide to oxygen. Despite this being general knowledge, deforestation continues at large with little consideration to its ramifications, which will result in the imminent extinction of a large number of plant species. According to earthsendangered.com, there are approximately 9000 endangered plant species worldwide, as of March 2015. This staggering figure necessitates establishing a plant species database dedicated to plant protection and data preservation.

Prior to the advances in computing technology, there exist a number of well-established plant classification methods. One of the earliest and most well-known methods is the Linnaeus taxonomy, whereby plants were subdivided into 24 orders of classes, which were used to identify but not represent natural groups of plants. At present, plant taxonomy still largely adopts the traditional classification methods.

Other modern methods of plant classification—molecular biology and morphological anatomy, to name a few—are deeply rooted in biology and chemistry.

The introduction of multidisciplinary studies such as image processing and machine learning in the computer science sector in the past decade has led to research in this field, producing non-manual plant classification systems. Compared to the abovementioned techniques, this method may prove superior, as leaf sampling and photographing said samples are less costly. The images of leaves can also be easily transferred to a computer so that the features can be extracted via image processing techniques.

The challenge, however, lies in the extraction and identification of discriminant features distinguishing plant features. While an image of a leaf will provide an array of data to be used with classification algorithms, not all the extracted features are relevant—in fact, too many irrelevant features inhibit the accuracy of results. Therefore, it is vital to identify relevant features in order to produce a classification model that can accurately classify a plant based on input vectors.

## **1.2 Problem statement**

As mentioned above, a lot of data, especially redundant data will cause the predicted output to deviate from the actual result. As such, feature selection is a must in the task of automatic plant classification in order to effectively eliminate irrelevant attributes, which minimizes the errors in classification.

In the past, multiple feature selection techniques have been known to work well, such as the PSO and ACO. However, the recent increase in the dimensionality of data poses a severe challenge to many existing feature selection techniques with respect to effectiveness and efficiency.

In the domain of agricultural robotics, automatically distinguishing between plant species is a challenging task, especially because some species are physically very similar to each other. The number of attributes one can use to classify a plants' species is thus too vast, whereby the use of an irrelevant/redundant attribute in the classification process may result in inaccurate results.

This study is done with the aim of identifying the most appropriate method of feature selection that applies to plant classification using SVM.

### 1.3 Objectives

- To study the efficiency of the Fast Correlated Based Filter(FCBF) feature selector algorithm in respect to predictive accuracy of the SVM in plant classification
- To propose the use of the Fast Correlated Based Filter as a feature selector for plant classification using SVM, if viable
- To compare the efficiency of the FCBF with other feature selecting algorithms.
- To evaluate the proposed techniques in feature selection for plant leaf classification

### 1.4 Scope

- Focus of project : plant leaf classification
- Experiments run using the Weka machine learning tool

- 3 feature selection algorithm (FCBF, CFS with PSO, and CFS with Best-First Search)
- The SMO classifier with Poly-kernel
- 2 sets of data:
  - Dataset 1 consists of 3 datasets, each with 1600 instances and 64 feature vectors: margin, texture and edge
  - Dataset 2 consists of 340 instances and 16 attributes

### **1.5 Project Significance**

- Pinpoint the suitability of proposed feature selection/preprocessing for the chosen datasets
- Qualitative analysis: guideline/reference for future work by other researchers

### **1.6 Expected results**

- Analysis of classification accuracy of SVM on features selected using abovementioned techniques
- Analysis of performance of feature selection technique

### **1.7 Conclusion**

The FCBF feature selector will be compared to two other feature selectors in the task of plant leaf classification in this project. The next chapter will cover literature reviews on this topic.

## **CHAPTER II**

### **LITERATURE REVIEW**

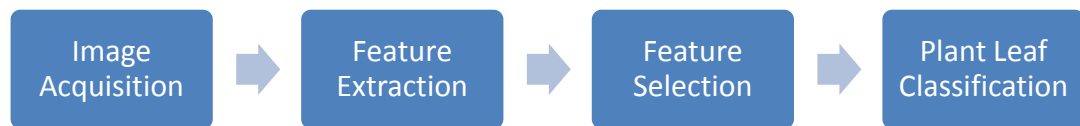
#### **2.1 Introduction**

Much research has been done on the application of Machine Learning in the agricultural field of plant leaf identification. While this chapter covers a variety of such application, this review focuses on the machine learning aspect of the process of plant leaf classification. The processes involved in plant leaf identification are: a) image acquisition, the acquisition of leaf images; b) Feature extraction: the extraction of data from leaf images; c) Feature selection: the selection of relevant, non-redundant features; d) Classification: the prediction of a plant species based on an array of attributes selected by the feature selector. Although this literature extensively covers the process of plant leaf identification, this paper will primarily focus on feature selection.

## 2.2 Plant Leaf Identification

The importance of an automated classification system for plants has been detailed in the previous chapter (chapter one).

For the past few years, researchers have worked on various techniques of automated plant classification. Generally, the process of classifying plants is as shown in Figure 2.1.



**Figure 2.1: Process of Plant Leaf Classification**

In *A leaf recognition algorithm for plant classification using probabilistic neural network* Wu et.al. (2007) proposed the probabilistic neural network (PNN) as a plant classification algorithm. In this work, 12 features were extracted and orthogonalized via Principle Component Analysis (PCA) to produce 5 principle components, which serve as inputs for the PNN. This method resulted in an accuracy of 90.312% when used on the Flavia dataset (which was provided by We et.al and available to the public). However, probabilistic neural networks are usually applied on smaller benchmarking datasets, as too large an array of inputs with which to train the network may result in overfitting. In *Neural networks for classification: a survey*, Zhang (2000) details the overfitting effect on test sets when using Probabilistic Neural Networks.

In *Svm-bdt pnn and fourier moment technique for classification of leaf shape*, Singh et.al. (2010) presented three methods of leaf classification: the SVM-BDT (Support Vector Machine utilizing Binary Decision Tree), PNN with PCA and Fourier moment, all of which utilized the same dataset as Wu et.al. (2007). In this research, the authors found that the SVM was superior to the PNN and Fourier moment due to its high generalization ability without the need of *a priori* knowledge.