# Tutorial

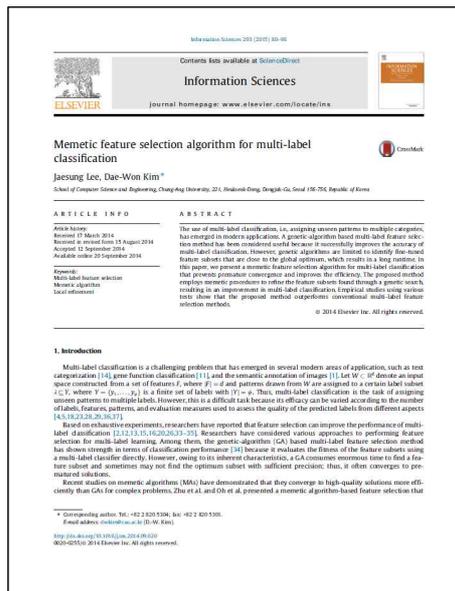## Memetic Feature Selection Algorithm for Multi-label Classification

Published from *Information Sciences* 293(1), 1 February 2015

### Jaesung Lee and Dae-Won Kim
### Department of Computer Science and Engineering, Chung-Ang University, Korea
9 November 2015

## Abstract

The use of multi-label classification, i.e., assigning unseen patterns to multiple categories, has emerged in modern applications. A genetic-algorithm based multi-label feature selection method has been considered useful because it successfully improves the accuracy of multi-label classification. However, genetic algorithms are limited to identify fine-tuned feature subsets that are close to the global optimum, which results in a long runtime. In this paper, we present a memetic feature selection algorithm for multi-label classification that prevents premature convergence and improves the efficiency. The proposed method employs memetic procedures to refine the feature subsets found through a genetic search, resulting in an improvement in multi-label classification. Empirical studies using various tests show that the proposed method outperforms conventional multi-label feature selection methods.

## Step 1. Download and unzip the file (programs.zip).



## Note

A. In this example, I unzipped the downloaded files to "**Z:\JSLee\programs**" folder.
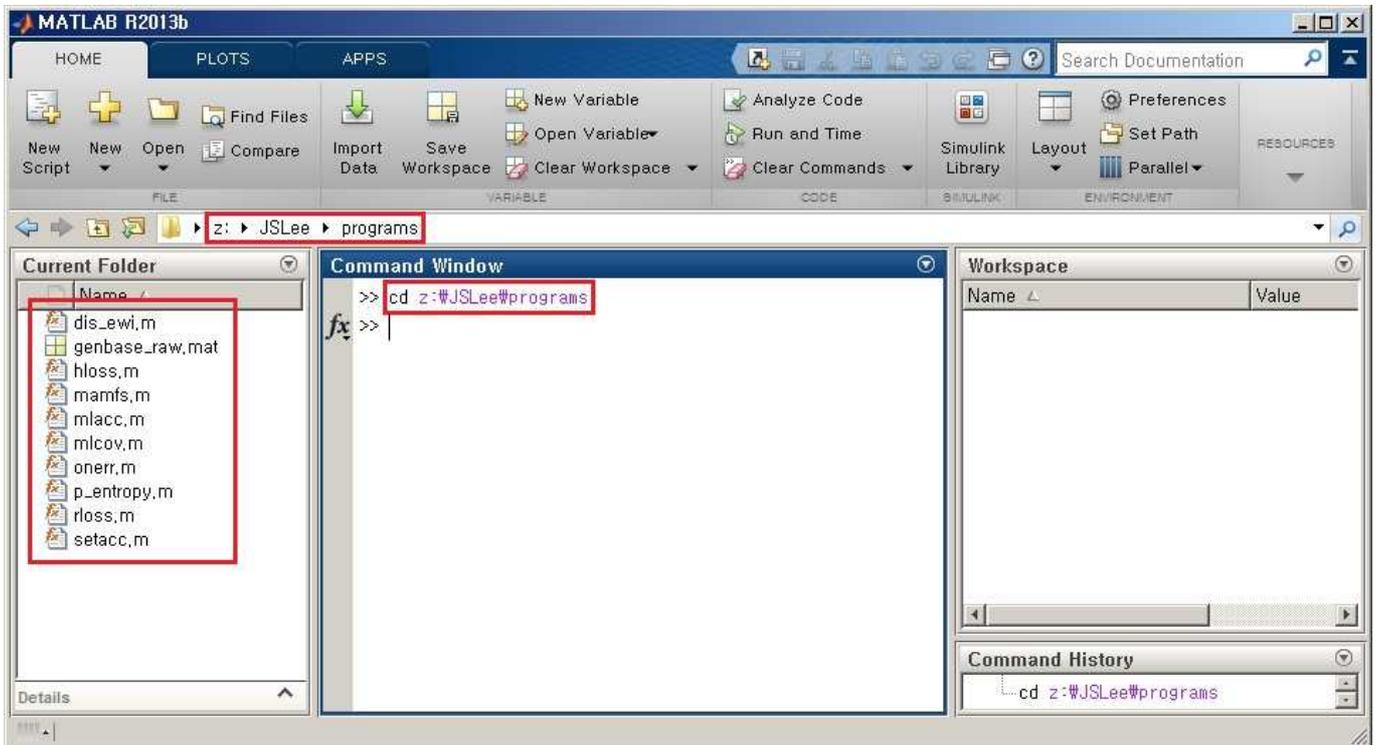
## Check points

A. "**programs.zip**" file contains 10 files in total.

B. "**genbase_raw.mat**" file contains an example data set that is publicly-available on the web site "http://mulan.sourceforge.net/download.html" with below references:

  * Mulan: Tsoumakas, G., Katakis, I., Vlahavas, I. (2010) "Mining Multi-label Data", Data Mining and Knowledge Discovery Handbook, O. Maimon, L. Rokach (Ed.), Springer, 2nd edition, 2010.

  * Genbase data set: S. Diplaris, G. Tsoumakas, P. Mitkas and I. Vlahavas. "Protein Classification with Multiple Algorithms," Proc. 10th Panhellenic Conference on Informatics (PCI 2005), pp. 448-456, Volos, Greece, November 2005.

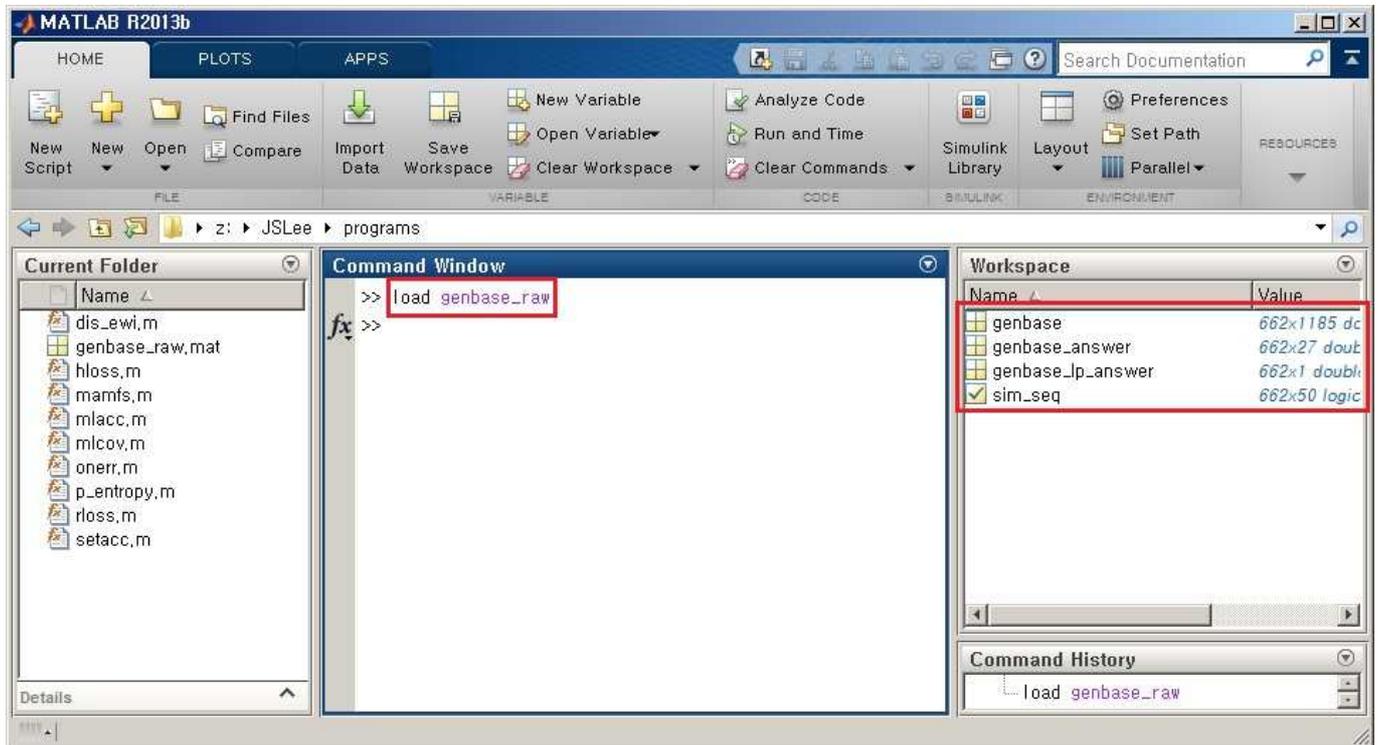## Step 2. Execute MATLAB and move to the work folder.



## Note

A. To move the current folder of MATLAB, type "**cd Z:\JSlee\programs**" to "**Command Window**".

## Check points

A. If you change the current folder correctly, you will see the unzipped files from "**Current Folder**" window.
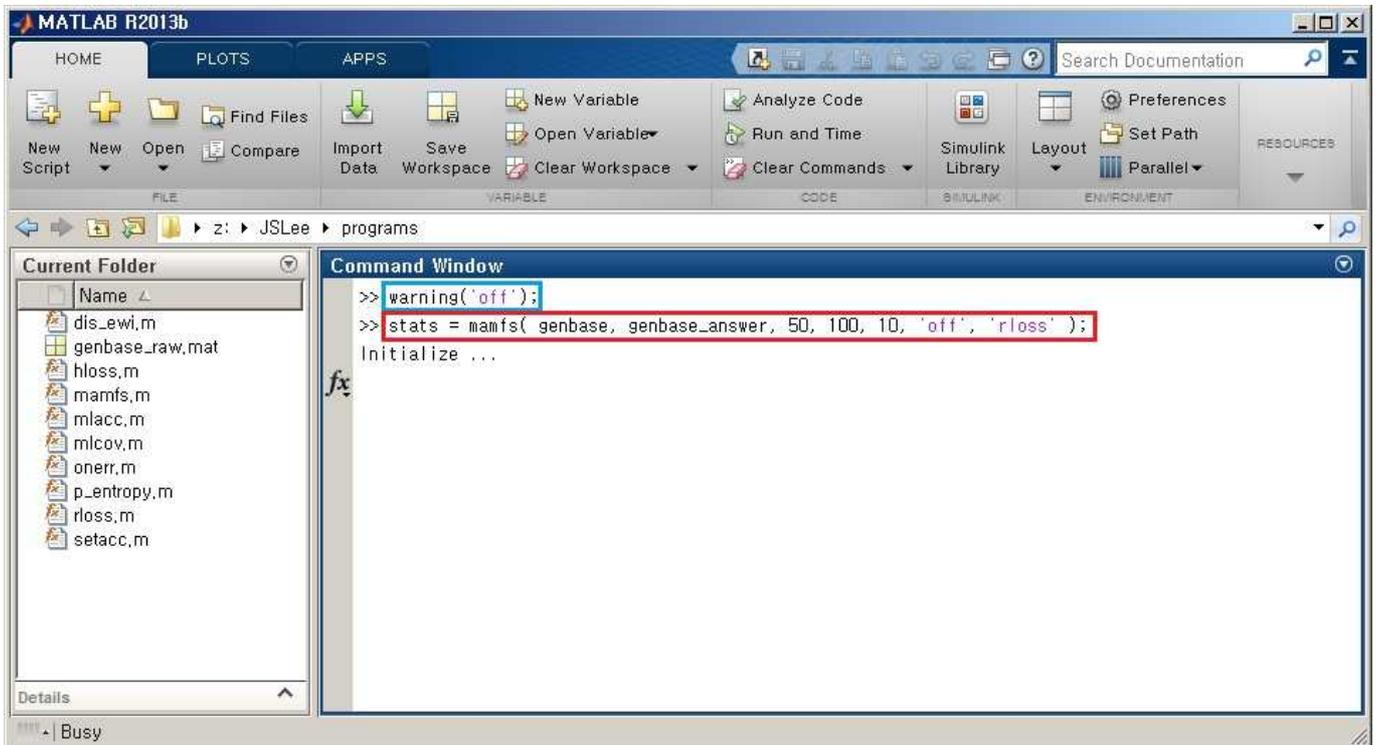
## Step 3. Load the example file (genbase_raw.m).



## Note

A. To load the example file, type "**load genbase_raw.m**" to "**Command Window**".

## Check points

A. If the example file is loaded successfully, you will see four variables (genbase, genbase_answer, genbase_lp_answer, and sim_seq) from "**Workspace**".

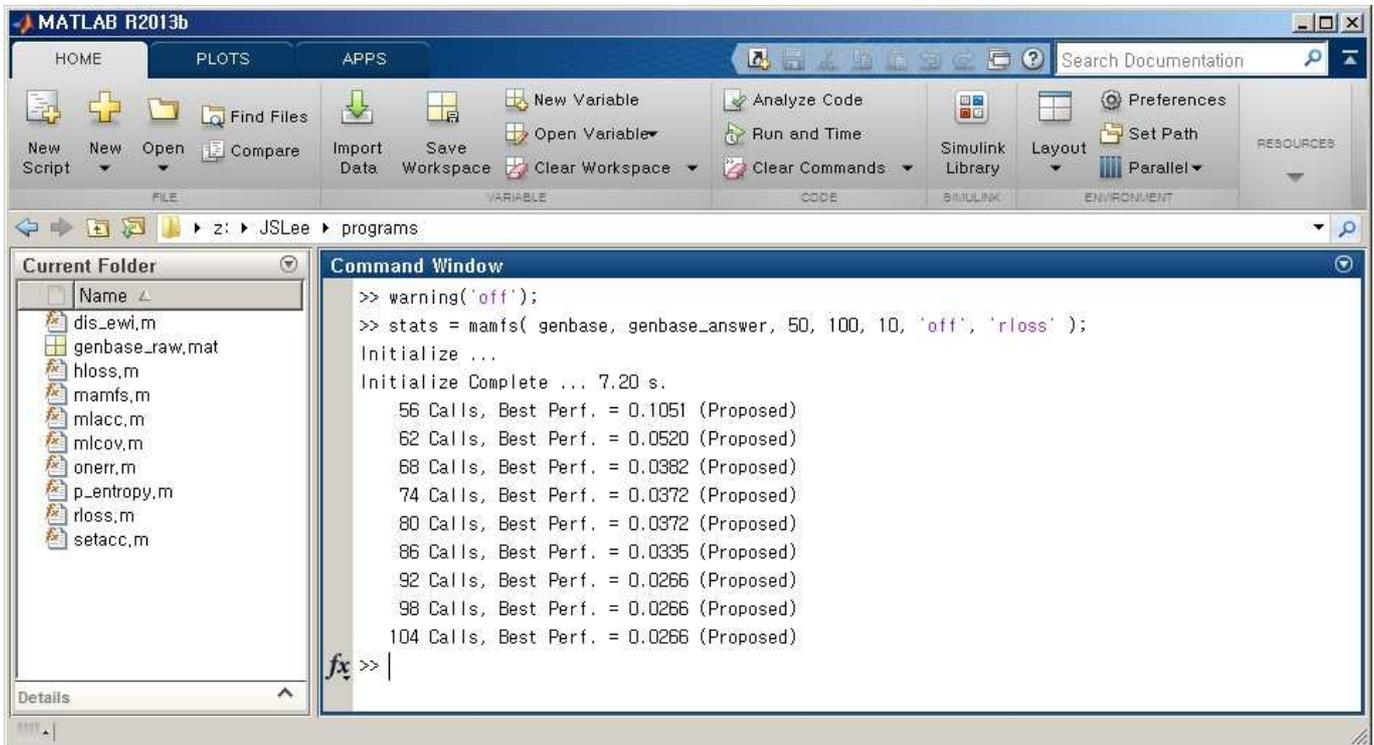# Step 4. Execute Memetic Algorithm for Multi-label Feature Selection (MAMFS).



## Note

A. "**warning('off');**" function turns off MATLAB global warning messages. For clarity, we recommend to run this option before actual execution of MAMFS.

B. To execute MAMFS, type "**stats = mamfs( genbase, genbase_answer, 50, 100, 10, 'off', 'rloss' );**" to "**Command Window**".

## Check points

A. If the MAMFS is run correctly, you will see "**Initialize ...**" message from "**Command Window**".

B. Detailed information about each variable is given below.

| Variable | Explanation |
|---|---|
| stats | 1×3 output cell matrix. Detailed information will be given in **Step 6**. |
| genbase | The data set matrix that is composed of 662 patterns and 1,185 features. |
| genbase_answer | The ground truth matrix that is composed of 27 labels (See **Step 3**). |
| 50 | The size of population (or the number of chromosomes in the population). |
| 100 | The maximum number of allowed fitness function calls. If MAMFS spends 100 fitness function calls during its execution, it will be terminated and returns "**stats**" variable. |
| 10 | The maximum size (or cardinality) of selected features. The size of feature subset will be smaller than specified value. |
| 'off' | If the input data matrix is a categorical (or binary) data set, it should be set to 'off'. In contrast, it should be set to 'on' if the input data matrix is a numerical data set. |
| 'rloss' | The name of fitness function. MAMFS allows six types of evaluation measures. <br> • 'hloss': Hamming loss    • 'rloss': Ranking loss    • 'mlacc': Multi-label accuracy <br> • 'setacc': Set accuracy    • 'onerr': One error    • 'mlcov': Coverage |

# Step 5. Wait a second until MAMFS returns the output variable.



## Note

A. Each line shows the fitness value of the best chromosome. For example, in this tutorial, MAMFS found a chromosome that gave 0.0266 of Ranking loss value after spending 104 fitness function calls (FFCs).
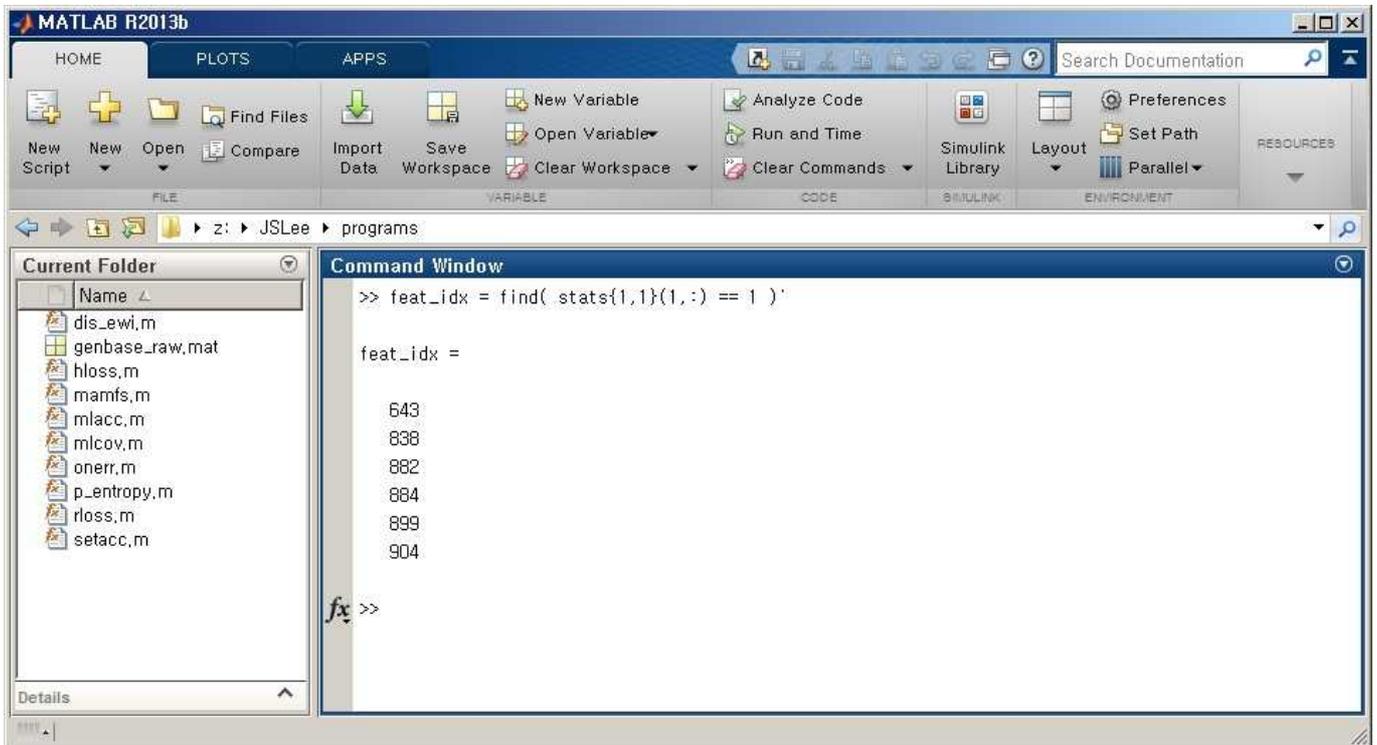
## Check points

A. MAMFS displays the best fitness value for each 6 FFCs. Specific fitness value for each FFCs can be easily obtained as below:

| Fitness Function Calls (FFCs) | Fitness (Ranking loss) |
|---|---|
| 56-61 | 0.1051 |
| 62-67 | 0.0520 |
| 68-73 | 0.0382 |
| 74-79 | 0.0372 |
| 80-85 | 0.0372 |
| 86-91 | 0.0335 |
| 92-97 | 0.0266 |
| 98-103 | 0.0266 |
| 104 (Final) | 0.0266 |

B. The gap of 6 FFCs comes from the genetic and local search process of MAMFS; 3 FFCs for local search and 3 FFCs for genetic search. During the creation or refinement process, MAMFS does not update the population, but spends FFCs for evaluating new chromosomes (or offspring).

# Step 6. Get the index of selected features.



## Note

A. After MAMFS returns "**stats**" variable, you can obtain the specific index of selected feature. To obtain the index, type "**feat_idx = find( stats{1,1}(1,:) == 1 )'**" to "**Command Window**".

B. In this tutorial, MAMFS returns a feature subset that is composed of 643th, 838th, 882th, 884th, 899th, 904th features in "**genbase**" data set.

## Check points

A. The "**stats**" variable is composed of three sub-cells. The meaning of each cell is given below:

| Position | Explanation |
|---|---|
| 1st | 50×1,185 matrix that represents the final population of MAMFS. It is sorted based on the given fitness value (Best to Worst). The values 0/1 represent selected/discarded features. |
| 2nd | The final fitness value of MAMFS. |
| 3rd | The 3rd cell is composed of five sub-cells. Meaning of each cell is given below: |

| | Position | Explanation |
|---|---|---|
| 3rd | 1st | The number of spent FFCs |
| | 2nd | The fitness value of corresponding FFCs described in the 1st cell. |
| | Remaining | Detailed fitness values correspond to given evaluation measure. If there is no specific statistics, it will be set to 'NaN' value. |