

# Using Genetic Algorithms for Data Mining Optimization in an Educational Web-based System

**GECCO 2003**

**Behrouz Minaei, William Punch**

Genetic Algorithm Research and Application Group  
Department of Computer Science and Engineering  
Michigan State University



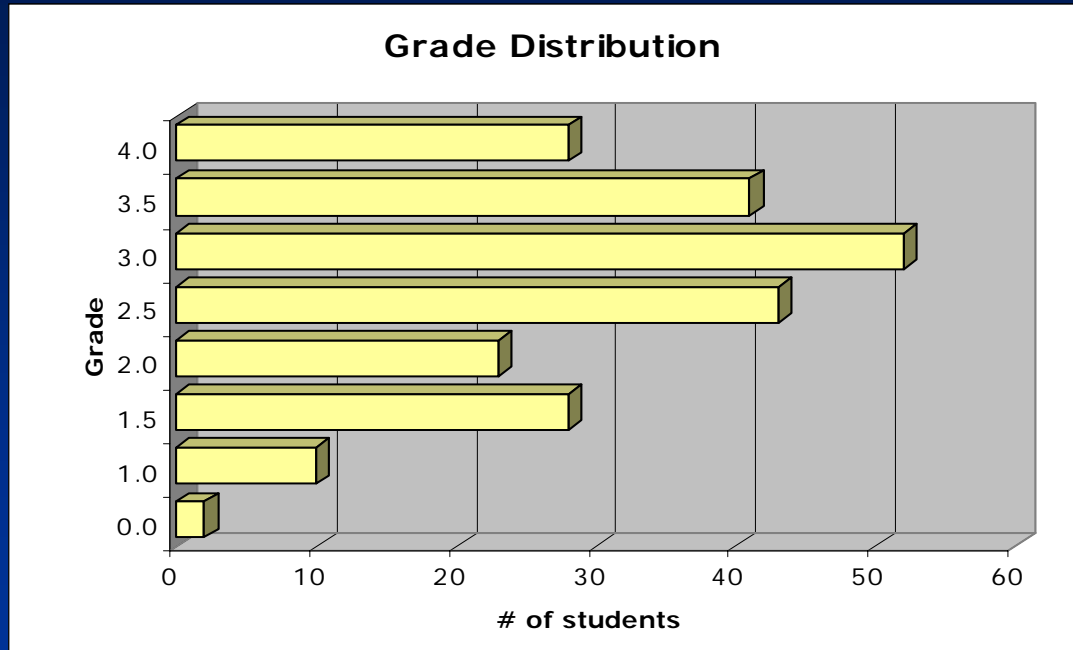
# Topics

- Problem Overview
- Classification Methods
  - Classifiers
  - Combination of Classifiers
- Weighting the features
- Using GA to choose best set of weights
- Experimental Results
- Conclusion & Next Steps

# Problem Overview

- This research is a part of the latest online educational system developed at Michigan State University (MSU), the *Learning Online Network with Computer-Assisted Personalized Approach (LON-CAPA)*.
- In LON-CAPA, we are involved with two kinds of large data sets:
  - Educational resources: web pages, demonstrations, simulations, individualized problems, quizzes, and examinations.
  - Information about users who create, modify, assess, or use these resources.
- Find *classes* of students. Groups of students use these online resources in a *similar* way.
- *Predict* for any individual student to which class he/she belongs.

# Data Set: PHY183 SS02



- 227 students
- 12 Homework sets
- 184 Problems

# Class Labels (3-ways)

## 2-Classes

1	<b>Passed</b>	Grade > 2.0	164	72.20%
2	<b>Failed</b>	Grade <= 2.0	63	27.80%

## 3-Classes

1	<b>High</b>	Grade >= 3.5	69	30.40%
2	<b>Middle</b>	2.0 < Grade < 3.5	95	41.80%
3	<b>Low</b>	Grade <= 2.0	63	27.80%

## 9-Classes

Class	1	2	3	4	5	6	7	8	9
Grade	0	0.5	1	1.5	2	2.5	3	3.5	4
# of students	2	0	10	28	23	43	52	41	28
Percentage	0.90%	0.00%	4.40%	12.40%	10.10%	18.90%	22.90%	18.00%	12.40%

# Extracted Features

1. Total number of correct answers.  
(Success rate)
2. Success at the first try
3. Number of attempts to get answer
4. Time spent until correct
5. Total time spent on the problem
6. Participating in the communication mechanisms

# Classifiers

- Non-Tree Classifiers (Using MATLAB)
  - Bayesian Classifier
  - 1NN
  - kNN
  - Multi-Layer Perceptron
  - Parzen Window
  - Combination of Multiple Classifiers (CMC)
  - *Genetic Algorithm (GA)*
- Decision Tree-Based Software
  - C5.0 (RuleQuest <<C4.5<<ID3)
  - CART (Salford-systems)
  - QUEST (Univ. of Wisconsin)
  - CRUISE [use an *unbiased* variable selection technique]

# GA Optimizer vs. Classifier

- Apply GA directly as a classifier
- Use GA as an optimization tool for resetting the parameters in other classifiers.
  - Most application of GA in pattern recognition applies GA as an optimizer for some parameters in the classification process.
  - Many researchers used GA in feature selection and feature extraction.
  - GA is applied to find an optimal set of feature weights that improve classification accuracy.
- Download GA Toolbox for MATLAB from:
  - <http://www.shef.ac.uk/~gaipp/ga-toolbox/>



# Fitness/Evaluation Function

- 5 classifiers:

1. Multi-Layer Perceptron

2 Minutes

2. Bayesian Classifier

3. 1NN

4. kNN

5. Parzen Window

- CMC

3 seconds

- Divide data into training and test sets (10-fold Cross-Val)

- **Fitness function: performance achieved by classifier**  
Error Rate in each round = 
$$\frac{\text{Total missclassified of test examples}}{\text{Total number of test examples}}$$

# Individual Representation

- The GA Toolbox supports binary, integer and floating-point chromosome representations.
- `Chrom = crtrp(NIND, FieldDR)` creates a random real-valued matrix of  $N \times d$ , `NIND` specifies the number of individuals in the population
- `FieldDR` is a matrix of size  $2 \times d$  and contains the boundaries of each variable of an individual.
- `FieldDR = [ 0 0 0 0 0 0 0; % lower bound`  
• `1 1 1 1 1 1 1]; % upper bound`
- `Chrom =` 0.23   0.17   0.95   0.38   0.06   0.26
- 0.35   0.09   0.43   0.64   0.20   0.54
- 0.50   0.10   0.09   0.65   0.68   0.46
- 0.21   0.29   0.89   0.48   0.63   0.89

# GA Parameters

- GGAP = 0.3
- XOVR = 0.7
- MUTR = 1/600
- MAXGEN = 500
- NIND = 1000
- SEL\_F = 'sus'
- XOVR\_F = 'recint'
- MUT\_F = 'mutbga'
- OBJ\_F = 'ObjFn2' / 'ObjFn3' / 'ObjFn9'

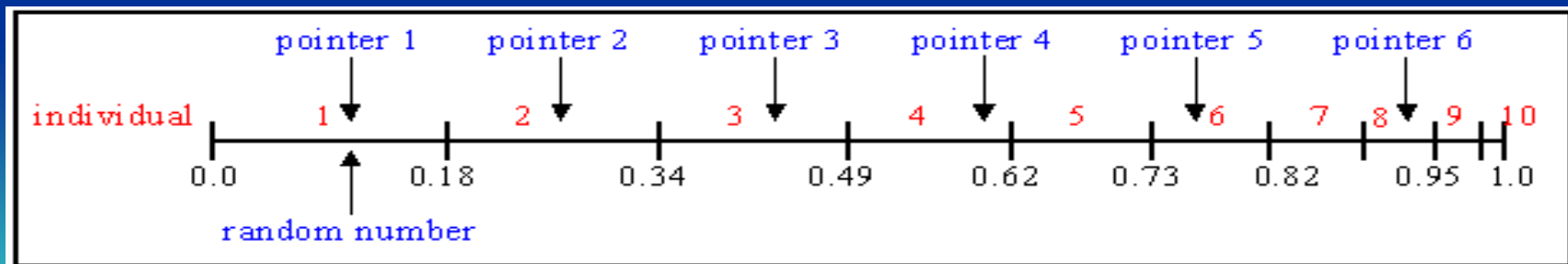
# Simple GA

```
% Create population
% Evaluate initial population,
gen = 0;
while gen < MAXGEN,
    % Assign fitness-value to entire population    FitnV = ranking(ObjV);
    % Select individuals for breeding SelCh = select(SEL_F, Chrom, FitnV,
    GGAP);
    % Recombine selected individuals (crossover)
    SelCh = recomb(XOV_F, SelCh,XOVR);
    % Perform mutation on offspring
    SelCh = mutate(MUT_F, SelCh, FieldD,
    MUTR);
    % Evaluate offspring, call objective function    ObjVSel = ObjFn2(data,
    SelCh);
    % Reinsert offspring into current population
    [Chrom
    ObjV]=reins(Chrom,SelCh,1,1,ObjV,ObjVSel);
    gen = gen+1;
    MUTR = MUTR+0.001;
end
```

# Selection - Stochastic Universal Sampling

- A form of stochastic universal sampling is implemented by obtaining a cumulative sum of the fitness vector,  $\text{FitnV}$ , and generating  $N$  equally spaced numbers between 0 and sum ( $\text{FitnV}$ ).
- Only one random number is generated, all the others used being equally spaced from that point.
- The index of the individuals selected is determined by comparing the generated numbers with the cumulative sum vector. The probability of an individual being selected is then given by

$$F(x_i) = \frac{f(x_i)}{\sum_{i=1}^{N_{ind}} f(x_i)}$$



# Crossover

OldChrom = [0.23 0.17 0.95 0.38 0.82 0.19; % parent1  
0.43 0.64 0.20 0.54 0.43 0.32] % parent2

Intermediate Recombination:

*NewChrom = recint(OldChrom)*

New values are produced by adding the scaled difference between the parent values to the first parent.

An internal table of scaling factors, Alpha, is produced, e.g.

Alpha = [0.13 0.50 0.32 0.16 0.23 0.06; % for offspring1  
0.12 0.54 0.44 0.26 0.27 0.10] % for offspring2

NewChrom = [ 0.40 0.92 0.86 0.33 0; % Alpha(1,:) parent1&2  
0.11 0.59 0.98 0.04]% Alpha(2,:) parent1&2

# Mutation

OldChrom = [0.23 0.81 0.17 0.66 0.28 0.30;  
0.43 0.96 0.25 0.64 0.10 0.23]

NewChrom = mutbga(OldChrom, FieldDR, [0.02 1.0]);

mutbga produces an internal mask table, MutMx = [0 0 0 1 0 1;  
0 0 1 0 0 0]

An second internal table, delta, specifies the normalized mutation  
step size,

delta = [0.02 0.02 0.02 0.02 0.02 0.02;  
0.01 0.01 0.01 0.01 0.01 0.01]

NewChrom = 0.23 0.81 0.17 0.66 0.28 0.30  
0.52 0.96 0.25 0.64 0.68 0.39

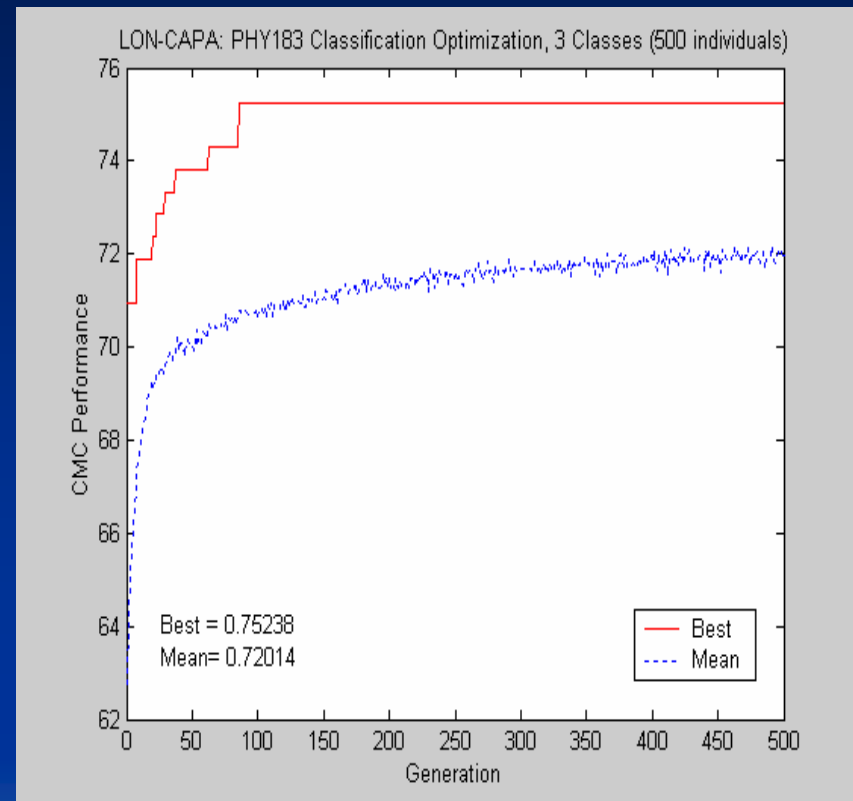
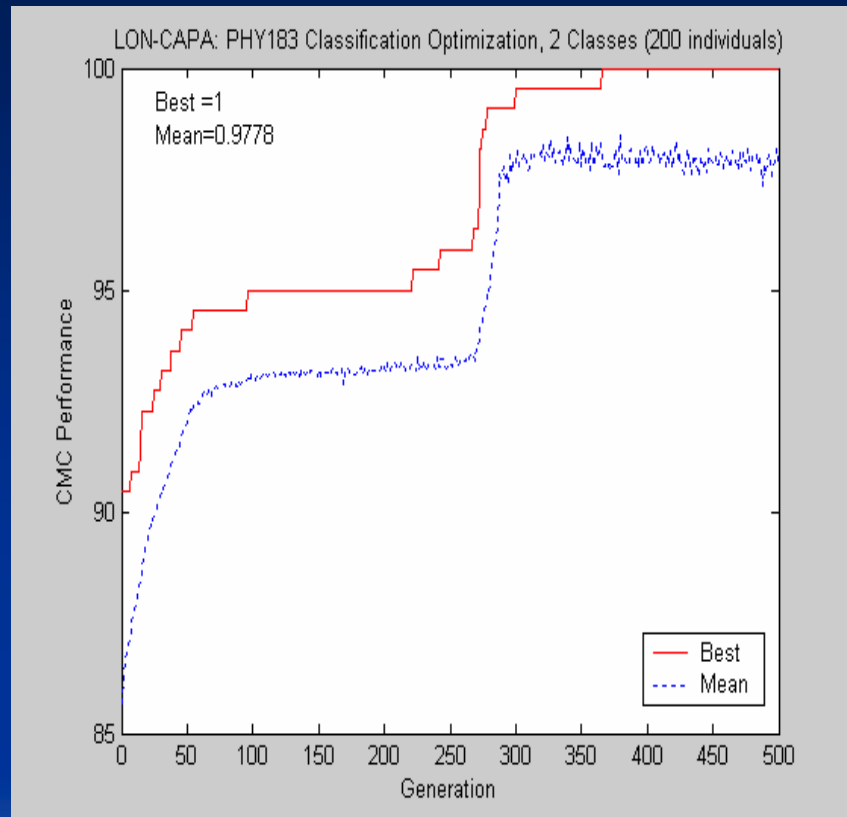
NewChrom - OldChrom shows the mutation steps = 0 0 0 4.61  
0 0 -7.51 -5.01

# Results without GA

		Performance %		
Classifier		2-Classes	3-Classes	9-Classes
Tree Classifier	<b>C5.0</b>	80.3	56.8	25.6
	<b>CART</b>	81.5	<b>59.9</b>	<b>33.1</b>
	<b>QUEST</b>	80.5	57.1	20.0
	<b>CRUISE</b>	81.0	54.9	22.9
Non-tree Classifier	<b>Bayes</b>	76.4	48.6	23.0
	<b>1NN</b>	76.8	50.5	29.0
	<b>kNN</b>	<b>82.3</b>	50.4	28.5
	<b>Parzen</b>	75.0	48.1	21.5
	<b>MLP</b>	79.5	50.9	-
	<b>CMC</b>	<b>86.8</b>	<b>70.9</b>	<b>51.0</b>



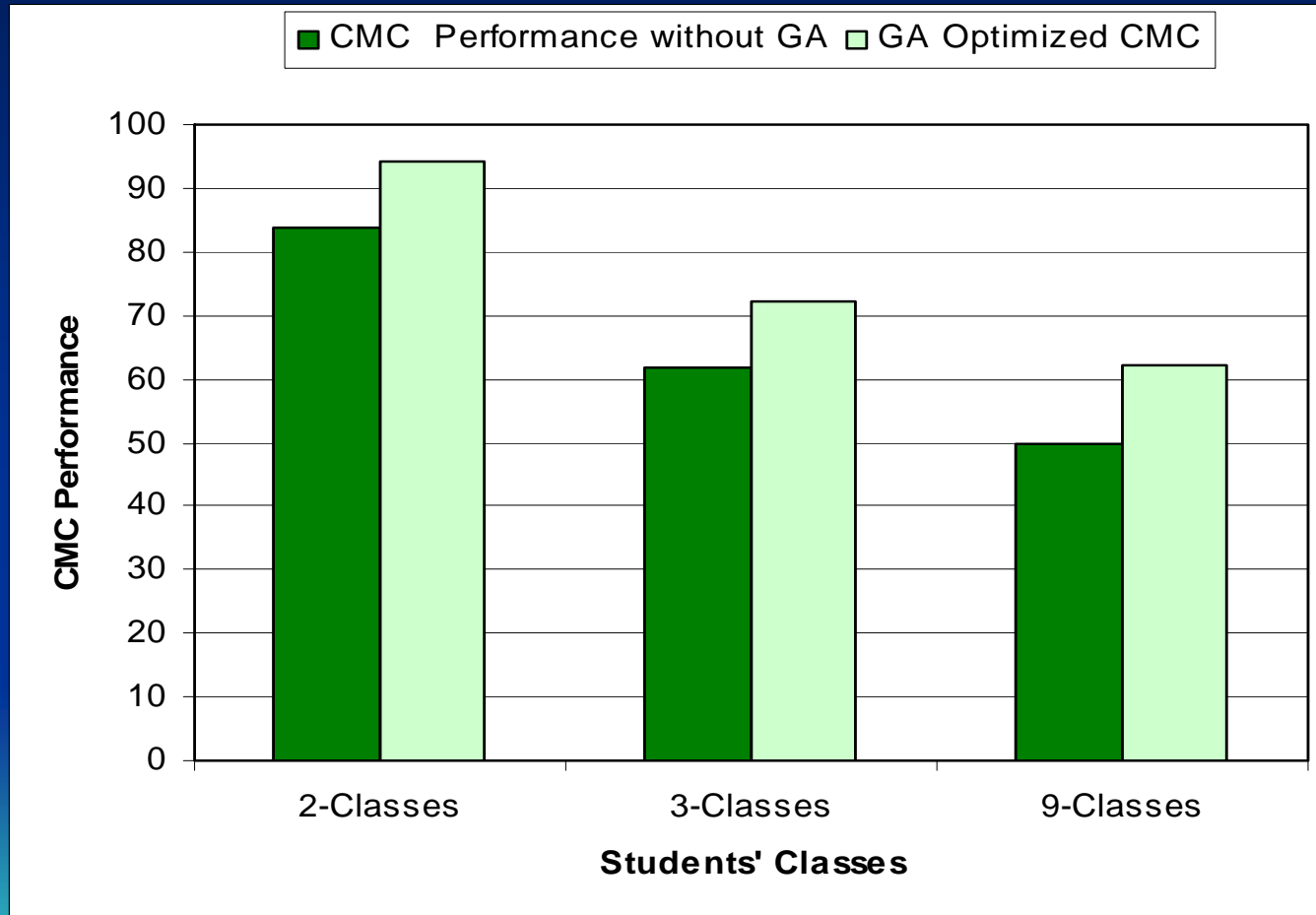
# Results of using GA



# GA Optimization Results

	Performance %		
Classifier	2- Classes	3- Classes	9- Classes
CMC of 4 Classifiers without GA	$83.87 \pm 1.73$	$61.86 \pm 2.16$	$49.74 \pm 1.86$
GA Optimized CMC, Mean individual	$94.09 \pm 2.84$	$72.13 \pm 0.39$	$62.25 \pm 0.63$
Improvement	$10.22 \pm 1.92$	$10.26 \pm 1.84$	$12.51 \pm 1.75$

# GA Optimization Results



# GA Optimization Results

Feature	Importance %
<u>Total Correct Answers</u>	100.00
<u>Total Number of Tries</u>	58.61
<u>First Got Correct</u>	27.70
<u>Time Spent to Solve</u>	24.60
<u>Total Time Spent</u>	24.47
<u>Communication</u>	9.21

# Summary

- Four classifiers used to segregate the students
- A combination of multiple classifiers leads to a significant accuracy improvement in all 3 cases.
- Weighted the features and used genetic algorithm to minimize the error rate.
- Genetic Algorithm could improve the prediction accuracy more than 10% in the case of 2 and 3-Classes and more than 12% in the case of 9-Classes.

# Next Steps

- Using Evolutionary Computation for extracting the features to find the optimal solution in LON-CAPA Data Mining
- Using Genetic programming to extract new features and improve prediction accuracy
- Apply EA to find the Frequency Dependency and Association Rules among the groups of problems  
*(Mathematical, Optional Response, Numerical, Java Applet, etc)*

# *Questions*