# COLUMBIA UNIVERSITY

MECE 4510 EVOLUTIONARY COMPUTATION AND DESIGN AUTOMATION

# Symbolic Regression

*Hanwen Zhao*
UNI: hz2547

supervised by
Dr. Hod LIPSON

Grace Hours Used: 53
Grace Hours Remaining: 43

October 22, 2018

# 1 Result Summary

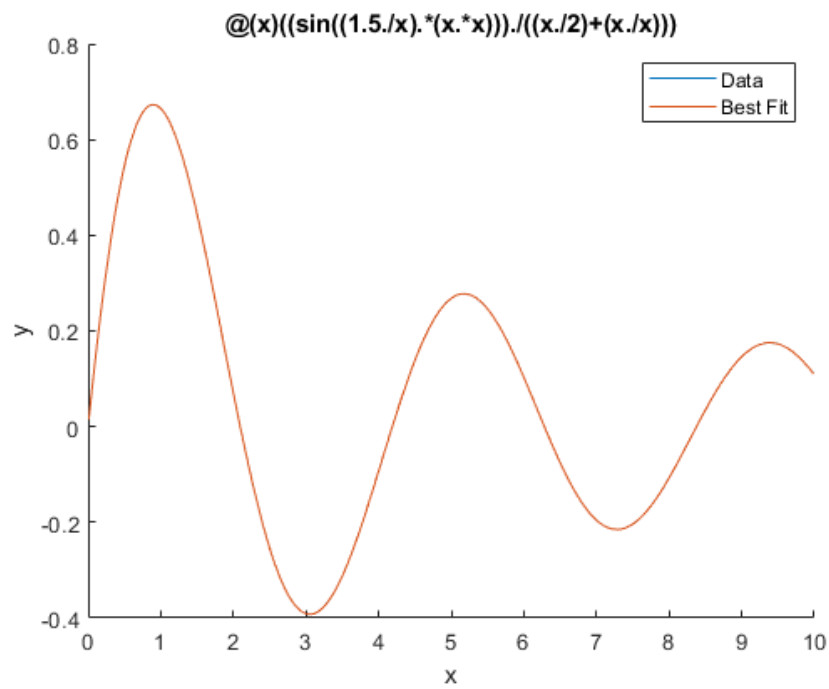| Result Table | | |
|---|---|---|
| Method | Evaluation Number | Best Error (MAE) |
| Random Search | 300000 | 0.0847 |
| Hill Climber | 300000 | 0.0863 |
| GP (Conventional Selection) | 14000 | 0.00002537 |
| GP (Conventional Selection) with Larger Population | 43000 | 0.00002537 |
| GP (Deterministic Crowding) | 129000 | 0.00002537 |



Figure 1: Performance Plot
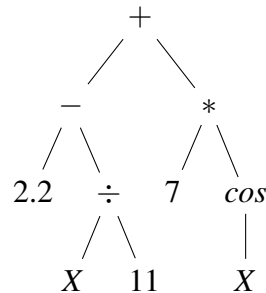
After simplification, the function can be expressed as:

$$f(x) = sin(1.5x)/(0.5x + 1) \tag{1}$$

# 2 Methods

For this homework, we have been asked to use genetic programming to perform symbolic regression. Our goal is to find the symbolic algebraic expression in the form of $y = f(x)$ that best fits a set of given 1000 $(x, y)$ points. Assume the symbolic regression only uses algebraic operators: $x$, $-$, $*$, $/$, *sine and cosine*, and real constants(in the range $\pm 10$), and the variable x.

## 2.1 Representation

The key part of genetic programming is to covert the program into a high level tree structure compare to genetic programming which the program was converted into simpler chromosome type. The tree structure is more powerful in terms of computer programming since trees can be easy evaluated in a recursive manner. For example, a function can be represented as a tree structure as following:



can represent the equation

$$(2.2 - \frac{x}{11}) + (7 * cos(x)) \tag{2}$$

For this assignment, we assume the maximum depth of the tree structure is less and equal than 5.

## 2.2 Random Search

Once we have our representation setup, the random search algorithm is straight forward. For each evaluation of the random search algorithm, a function tree is generated with random depth from two to five. The random search is the baseline for performance comparison between the hill climber algorithm and our genetic programmings.

## 2.3 Random Hill Climber

A Random Hill Climber is basically a random search with simple decision making capability. Between each evaluation, the mutation process is applied to the function. During the mutation process, it will randomly select a valid mutation point, generate a new sub-tree and replace into new function tree. In my implementation, the probability decide whether it always go to a better solution or not. Sometime go to a worse solution can avoid hill climber stuck at the local maximum.

## 2.4 Genetic Programming with Variations

Similar to Evolutionary Algorithm, we can implement following types of operators into the Genetic Programming: selection, crossover, and mutation. For variation, Deterministic Crowding and different size of population was applied during my implementation. Here are some details about the techniques I used in my implementation:

- Selection: For conventional selection Genetic Programming, during each generation, a default of 50% parents will be selected to generate 50 offsprings.

- Crossover: The crossover was applied into both conventional selection method and deterministic crowding method. The algorithm will first choose two valid crossover node from each parent and switch the subtrees between two parents in order to generate two offSprings.

- Mutation: During the mutation process, the algorithm will first randomly pick a valid node for mutation, then generate a new tree based on the depth of the mutation node.

- Deterministic Crowding: In order to maintain the diversity for genetic programming, we need some methods to maintain useful diversity for genetic programming to work better. Crowding is one of the popular method, it only replace individuals that are similar. More specific, deterministic crowding compare the similarity between two parents and two offsrpings, replace the one has higher similarity.

---
**Algorithm 1** Deterministic Crowding

---
1: **procedure** MYPROCEDURE
2:      **if** $d(p_1, c_1) + d(p_2, c_2) < d(p_1, c_2) + d(p_2, c_1)$ **then**
3:          *compare $c_1$ to $p_1$ and $c_2$ to $p_2$ and replace parents if offspring better*
4:      **else**
5:          *compare $c_1$ to $p_2$ and $c_2$ to $p_1$ and replace parents if offspring better*

---

## 2.5 Analysis of Performance

Overall, all GPs performed very well since they all can find the optimal solution. However, the random search and hill climber did not perform as expected. The main reason for the failure of hill climber because it was badly guided. Like random search, hill climber would perform well if it is lucky. On the other hand, the deterministic crowding for genetic programming is extremely powerful. From both performance plot as well as the diversity plot, the crowding helped to maintain diversity in a better manner compare to conventional selection. Therefore, on the performance plot, the crowding starts slow than conventional selection, but it could find better solution eventually.

3

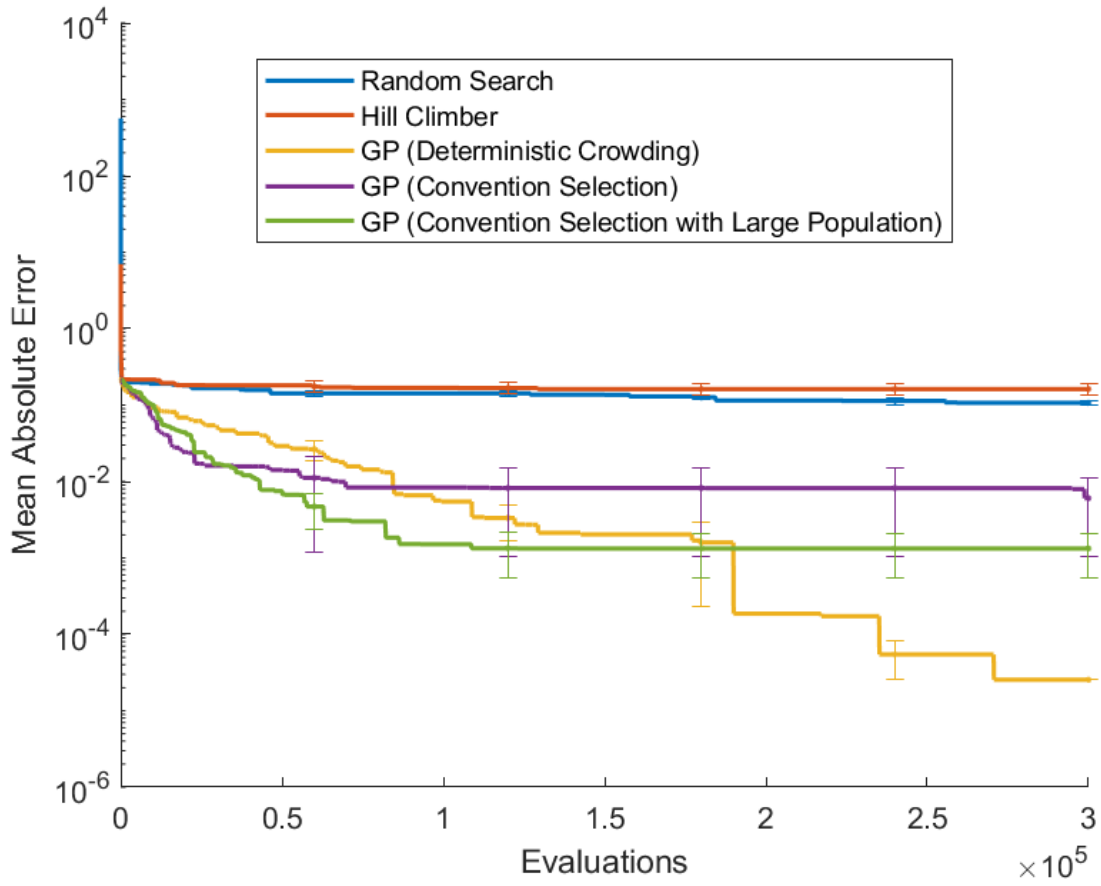# 3 Performance

## 3.1 Performance Plots



Figure 2: Performance Plot
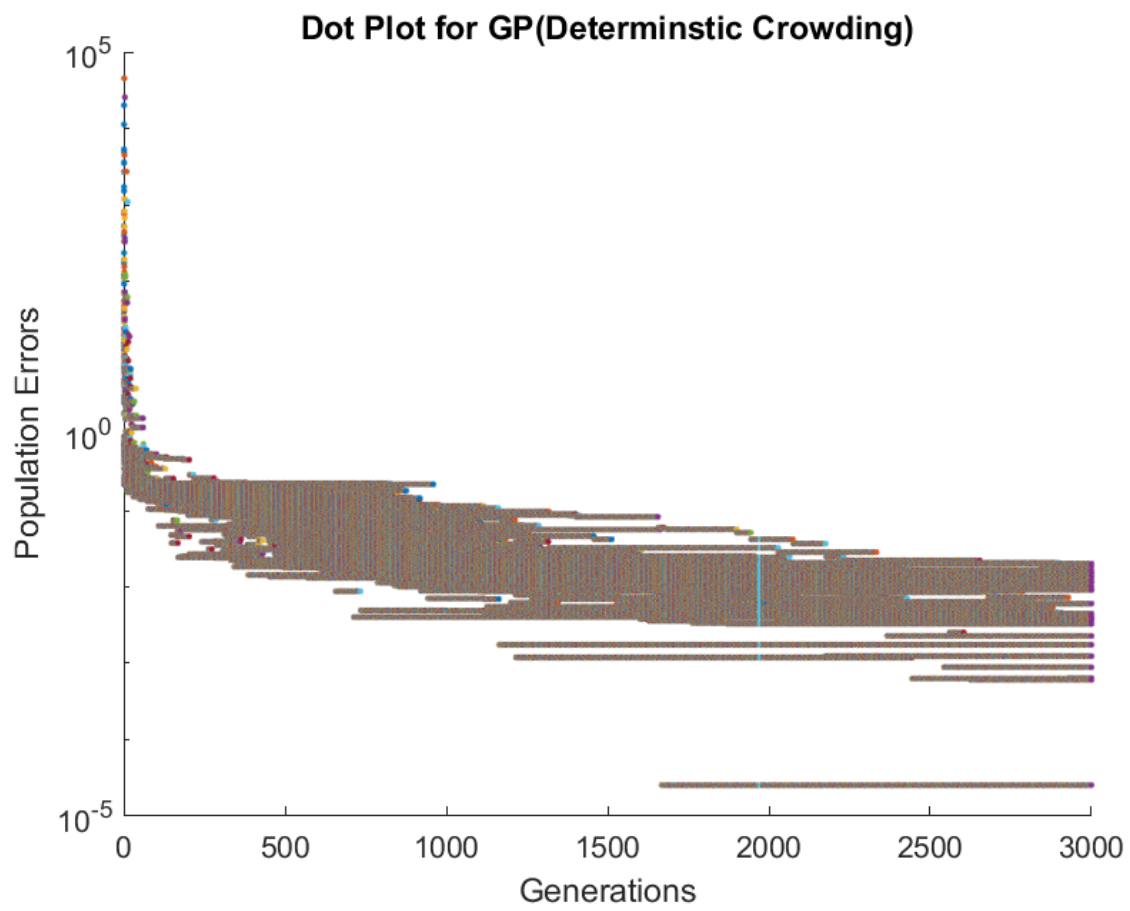
## 3.2   Dot Plot



Figure 3: Dot Plot

## 3.3 Diversity Plot
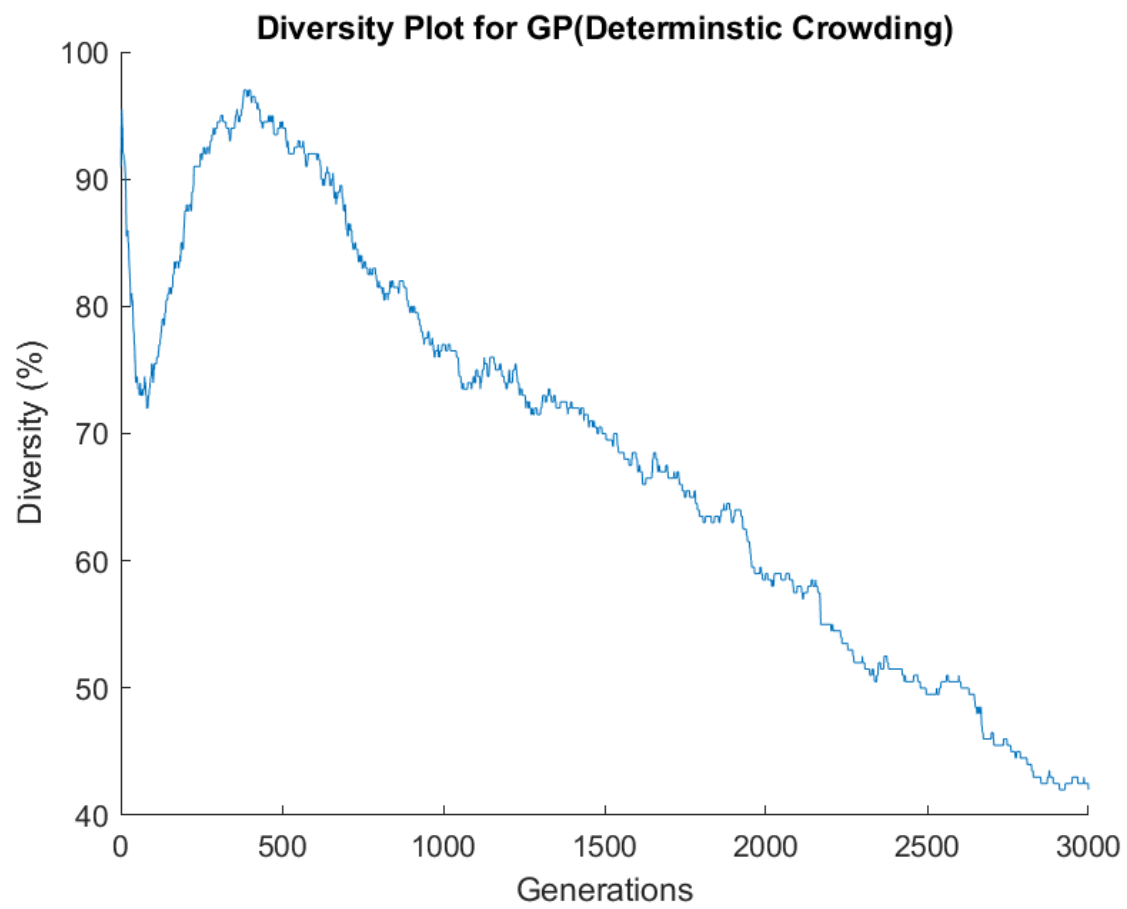


Figure 4: Diversity Plot

## 3.4 Convergence Plot



Figure 5: Convergence Plot

## 3.5 Simpler Problem Tested

During the debugging process of Genetic Programming, simpler problem was tested to ensure the GP can run without any bug. For example, a test data set of $sin(x)$ was used. In addition, all GPs can find the optimal solution within 10 generations.

Figure 6: Simple Problem Test

## 3.6   Validation

The validation plot did not show much difference between training data and testing data since the data we were given has negligible noise.



Figure 7: Validation Plot

# 4   Appendix

```matlab
1  function [] = HW2()
2  operator = {'+','-','.*','./','sin','cos'};
3  varConst = {'c','x'};
4  maxLevel = 4;
5  tic
6  rng('shuffle')
7  data = csvread('function1.csv');
8  %trainingIndex = sort(randperm(1000, 100));
9  trainingIndex = 1:500;
10 trainingData = data(trainingIndex,:);
11 validationData = data(setdiff(1:length(data),trainingIndex),:);
12 randomSearch(operator, varConst, maxLevel, trainingData, 300000, 4);
13 hillClimber(operator, varConst, maxLevel, trainingData, 300000, 4);
14 GP1(operator, varConst, maxLevel, trainingData, validationData, 200, 3000, 4, 'true', 'dot');
15 GP2(operator, varConst, maxLevel, trainingData, 200, 3000, 4, 'true');
16 GP2_LP(operator, varConst, maxLevel, trainingData, 400, 1500, 4, 'true');
17 toc
18 end
19
20 function [] = GP2_LP(operator, varConst, maxLevel, data, popSize, n, repeat, print)
21 parfor r = 1:repeat
22     if strcmp(print,'true')
23         fileID = fopen(strcat('GP2_LP_',int2str(r),'.txt'),'wt');
24     end
25     population = cell([popSize,2^maxLevel]);
26     for ii = 1:popSize
27         population(ii,1:end-1) = heapGeneration(operator, varConst, maxLevel);
28     end
29     popError = zeros(popSize,1);
30     bestError = inf;
31     for i = 1:n
32         for ii = 1:popSize
33             population{ii,2^maxLevel} = MAE_Cal(population(ii,1:end-1),data);
34             popError(ii) = population{ii,2^maxLevel};
35             %heapString(population(ii,1:end-1))
36         end
37         currentError = min(popError);
38         population = sortrows(population,2^maxLevel);
39         selectPop = population(1:popSize/2,1:end-1);
40         offSpringPop = cell([popSize/2,2^maxLevel-1]);
41         bestError = min(currentError,bestError);
42         if strcmp(print,'true')
43             fprintf(fileID, '%d %10.8f %10.8f \n', i*popSize/2, currentError, bestError);
44         end
45         for jj = 1:popSize/4
46             parent1 = selectPop(randi(length(selectPop)),:);
47             parent2 = selectPop(randi(length(selectPop)),:);
48             [offSpring1, offSpring2] = crossOver(parent1, parent2, maxLevel);
49             offSpring1 = mutation(offSpring1, operator, varConst, maxLevel);
50             offSpring2 = mutation(offSpring2, operator, varConst, maxLevel);
51             offSpringPop(jj*2-1:jj*2,:) = [offSpring1;offSpring2];
52         end
53         population(:,1:end-1) = [selectPop;offSpringPop];
54     end
55     population = population(:,1:end-1);
56     for i = 1:popSize
57         popError(i) = MAE_Cal(population(i,:),data);
58     end
59 %     bestHeap = population(1,:);
60 %     bestHeapString = heapString(bestHeap);
61 %     figure
62 %     hold on
63 %     plot(data(:,1),data(:,2))
64 %     plot(data(:,1),evalHeap(bestHeap, data))
65 %     legend('Data','Best Fit')
```
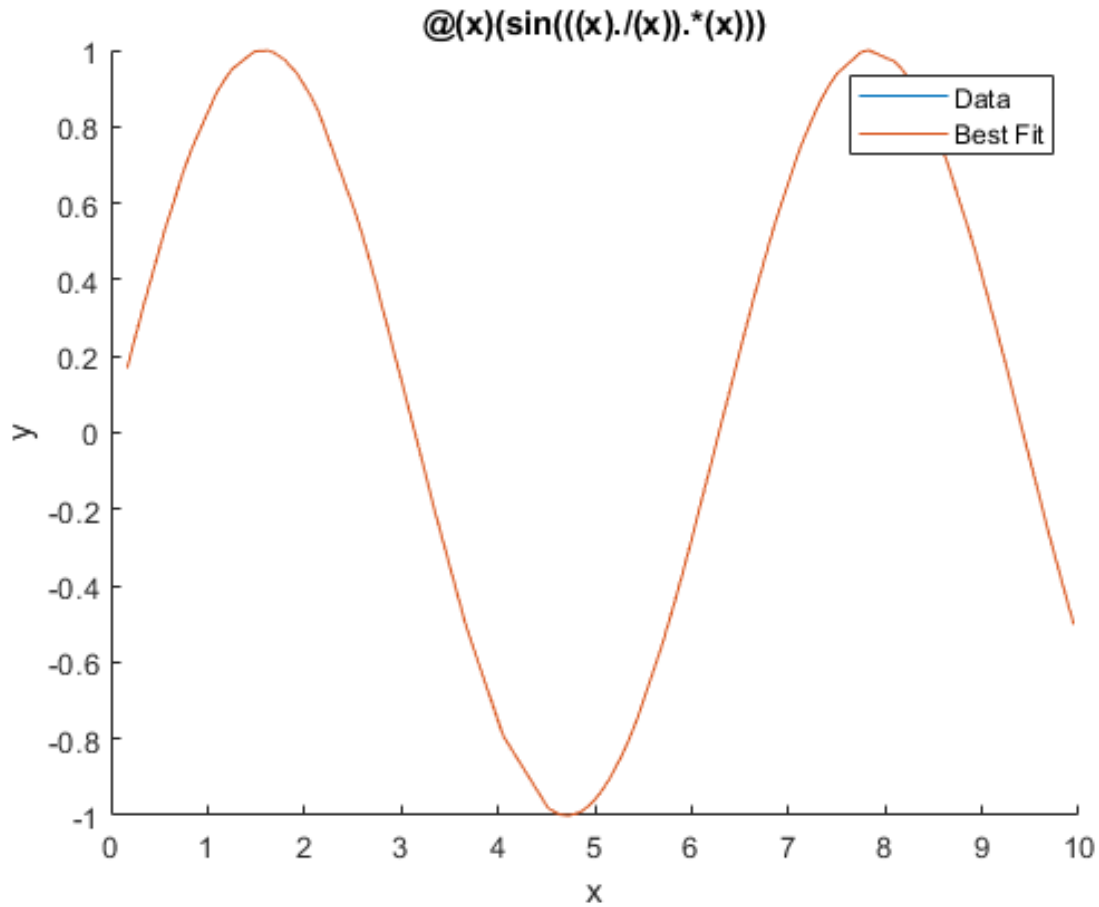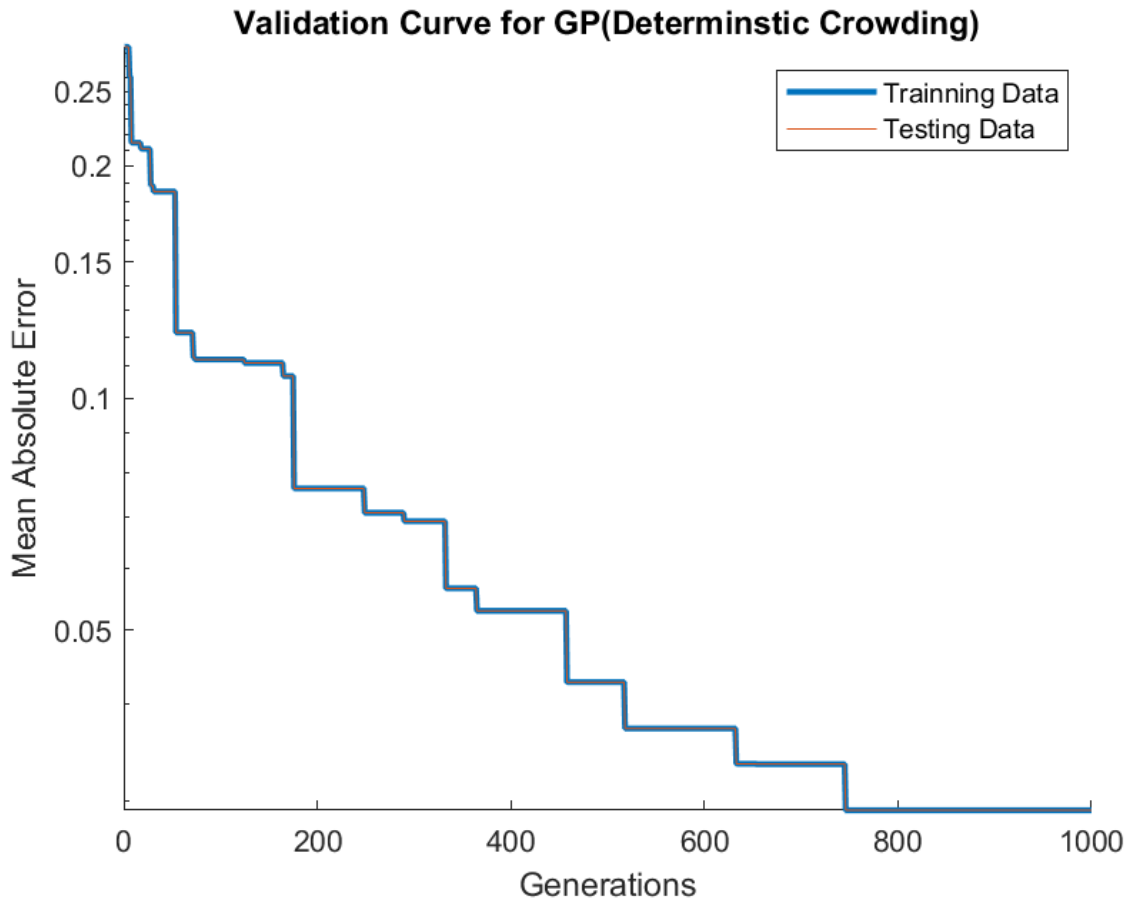
10

```matlab
66  %        xlabel('x')
67  %        ylabel('y')
68  %        title(bestHeapString)
69         index = find(popError == min(popError));
70         fprintf("GP2: Eval: %d, Error: %10.8f.\n", n, popError(index(1)))
71  end
72  end
73
74  function [] = GP2(operator, varConst, maxLevel, data, popSize, n, repeat, print)
75  parfor r = 1:repeat
76      if strcmp(print,'true')
77          fileID = fopen(strcat('GP2_',int2str(r),'.txt'),'wt');
78      end
79      population = cell([popSize,2^maxLevel]);
80      for ii = 1:popSize
81          population(ii,1:end-1) = heapGeneration(operator, varConst, maxLevel);
82      end
83      popError = zeros(popSize,1);
84      bestError = inf;
85      for i = 1:n
86          for ii = 1:popSize
87              population{ii,2^maxLevel} = MAE_Cal(population(ii,1:end-1),data);
88              popError(ii) = population{ii,2^maxLevel};
89              %heapString(population(ii,1:end-1))
90          end
91          currentError = min(popError);
92          population = sortrows(population,2^maxLevel);
93          selectPop = population(1:popSize/2,1:end-1);
94          offSpringPop = cell([popSize/2,2^maxLevel-1]);
95          bestError = min(currentError,bestError);
96          if strcmp(print,'true')
97              fprintf(fileID, '%d %10.8f %10.8f \n', i*popSize/2, currentError, bestError);
98          end
99          for jj = 1:popSize/4
100             parent1 = selectPop(randi(length(selectPop)),:);
101             parent2 = selectPop(randi(length(selectPop)),:);
102             [offSpring1, offSpring2] = crossOver(parent1, parent2, maxLevel);
103             offSpring1 = mutation(offSpring1, operator, varConst, maxLevel);
104             offSpring2 = mutation(offSpring2, operator, varConst, maxLevel);
105             offSpringPop(jj*2-1:jj*2,:) = [offSpring1;offSpring2];
106         end
107         population(:,1:end-1) = [selectPop;offSpringPop];
108     end
109     population = population(:,1:end-1);
110     for i = 1:popSize
111         popError(i) = MAE_Cal(population(i,:),data);
112     end
113 %    bestHeap = population(1,:);
114 %    bestHeapString = heapString(bestHeap);
115 %    figure
116 %    hold on
117 %    plot(data(:,1),data(:,2))
118 %    plot(data(:,1),evalHeap(bestHeap, data))
119 %    legend('Data','Best Fit')
120 %    xlabel('x')
121 %    ylabel('y')
122 %    title(bestHeapString)
123     index = find(popError == min(popError));
124     fprintf("GP2: Eval: %d, Error: %10.8f.\n", n, popError(index(1)))
125 end
126 end
127
128 function [] = GP1(operator, varConst, maxLevel, data, validationData, popSize, n, repeat, print, ...
        type)
129 dotData = zeros(n,popSize);
130 errorAndValidation = zeros(n,2);
131 for r = 1:repeat
132     if strcmp(print,'true')
```

```matlab
133            fileID = fopen(strcat('GP1_',int2str(r),'.txt'),'wt');
134        end
135        population = cell([popSize,2^maxLevel-1]);
136        for ii = 1:popSize
137            population(ii,:) = heapGeneration(operator, varConst, maxLevel);
138        end
139        popError = zeros(popSize,1);
140        validationError = zeros(popSize,1);
141        bestError = inf;
142        for i = 1:n
143            for ii = 1:popSize
144                popError(ii) = MAE_Cal(population(ii,:),data);
145                validationError(ii) = MAE_Cal(population(ii,:),validationData);
146                %heapString(population(ii,:))
147            end
148            mean(popError,'omitnan');
149            currentError = sort(popError);
150            validationError = sort(validationError);
151            bestError = min(currentError(1),bestError);
152            errorAndValidation(i,:) = [currentError(1) validationError(1)];
153            if strcmp(print,'true')
154                fprintf(fileID, '%d %10.8f %10.8f \n', i*popSize/2, currentError(1), bestError);
155            end
156            if strcmp(type,'dot')
157                dotData(i,:) = popError;
158            end
159            % random select two parents
160            for j = 1:popSize/2
161                index = randperm(popSize,2);
162                parent1 = population(index(1),:);
163                parent2 = population(index(2),:);
164                % crossover to get two offsprings
165                [offSpring1, offSpring2] = crossOver(parent1, parent2, maxLevel);
166                % perfome mutation
167                offSpring1 = mutation(offSpring1, operator, varConst, maxLevel);
168                offSpring2 = mutation(offSpring2, operator, varConst, maxLevel);
169                % perfome deterministic crowding
170                dp1c1 = similarCal(parent1, offSpring1, data);
171                dp2c2 = similarCal(parent2, offSpring2, data);
172                dp1c2 = similarCal(parent1, offSpring2, data);
173                dp2c1 = similarCal(parent2, offSpring1, data);
174                c1 = MAE_Cal(offSpring1, data); p1 = MAE_Cal(parent1, data);
175                c2 = MAE_Cal(offSpring2, data); p2 = MAE_Cal(parent2, data);
176                if dp1c1+dp2c2 < dp1c2+dp2c1
177                    if c1 < p1
178                        population(index(1),:) = offSpring1;
179                    end
180                    if c2 < p2
181                        population(index(2),:) = offSpring2;
182                    end
183                else
184                    if c1 < p2
185                        population(index(2),:) = offSpring1;
186                    end
187                    if c2 < p1
188                        population(index(1),:) = offSpring2;
189                    end
190                end
191            end
192        end
193        for i = 1:popSize
194            popError(i) = MAE_Cal(population(i,:),data);
195        end
196    %     sortError = sort(popError);
197    %     [~,errorIndex] = ismember(sortError,popError);
198    %     sortPop = population(errorIndex,:);
199    %     bestHeap = sortPop(1,:);
200    %     bestHeapString = heapString(bestHeap);
```

```matlab
201  %      figure
202  %      hold on
203  %      plot(data(:,1),data(:,2))
204  %      plot(data(:,1),evalHeap(bestHeap, data))
205  %      legend('Data','Best Fit')
206  %      xlabel('x')
207  %      ylabel('y')
208  %      title(bestHeapString)
209      index = find(popError == min(popError));
210      fprintf("GP1: Eval: %d, Error: %10.8f.\n", n, popError(index(1)))
211      save('dotData','dotData')
212      save('errorAndValidation2','errorAndValidation')
213  end
214  end
215
216  function sim = similarCal(heap1, heap2, data)
217  sim = abs(MAE_Cal(heap1,data)-MAE_Cal(heap2,data));
218  end
219
220  function [offSpring1, offSpring2] = crossOver(parent1, parent2, maxLevel)
221  maxCrossOverLimit = min(max(find(~cellfun(@isempty,parent1))),max(find(~cellfun(@isempty,parent2)
         ))));
222  offSpring1 = parent1;
223  offSpring2 = parent2;
224  maxCrossOverLevel = floor(log2(maxCrossOverLimit))+1;
225  while 1
226  crossOverLevel = randi(maxCrossOverLevel);
227  selection = 2^(crossOverLevel-1):2^(crossOverLevel)-1;
228  crossOverPoint1 = selection(randi(length(selection)));
229  crossOverPoint2 = selection(randi(length(selection)));
230  if ~isempty(parent1{crossOverPoint1}) && ~isempty(parent2{crossOverPoint2})
231      if isnumeric(parent1{crossOverPoint1}) && isnumeric(parent2{crossOverPoint2})
232          break
233      elseif strcmp(parent1{crossOverPoint1},'x') && strcmp(parent2{crossOverPoint2},'x')
234          break
235      elseif isnumeric(parent1{crossOverPoint1}) && strcmp(parent2{crossOverPoint2},'x')
236          break
237      elseif isnumeric(parent2{crossOverPoint2}) && strcmp(parent1{crossOverPoint1},'x')
238          break
239      elseif ischar(parent1{crossOverPoint1}) && ischar(parent2{crossOverPoint2})
240          break
241      end
242  end
243  end
244  crossOverLocations1 = searchChildren(crossOverPoint1, maxLevel);
245  crossOverLocations2 = searchChildren(crossOverPoint2, maxLevel);
246  subHeap1 = parent1(crossOverLocations1);
247  subHeap2 = parent2(crossOverLocations2);
248  for i = 1:length(crossOverLocations1)
249      offSpring1(crossOverLocations1(i)) = subHeap2(i);
250      offSpring2(crossOverLocations2(i)) = subHeap1(i);
251  end
252  end
253
254  function [] = hillClimber(operator, varConst, maxLevel, data, n, repeat)
255  parfor r = 1:repeat
256      fileID = fopen(strcat('RMHC_',int2str(r),'.txt'),'wt');
257      heap = heapGeneration(operator, varConst, maxLevel);
258      oldError = MAE_Cal(heap,data);
259      bestError = inf;
260      for i = 1:n
261          newHeap = mutation(heap, operator, varConst, maxLevel);
262          newError = MAE_Cal(newHeap, data);
263          if newError < oldError
264              heap = newHeap;
265              oldError = newError;
266              bestError = newError;
267          end
```

```matlab
268            fprintf(fileID, '%d %10.8f %10.8f \n', i, newError, bestError);
269        end
270        fprintf("RMHC: Eval: %d, Error: %10.8f.\n", n, bestError)
271    end
272    end
273
274    function heap = mutation(heap, operator, varConst, maxLevel)
275    % the mutation should be able happen non-empty node
276    mutateIndex = find(~cellfun(@isempty,heap));
277    mutationPoint = mutateIndex(randi(length(mutateIndex)-1)+1);
278    mutateLocations = searchChildren(mutationPoint, maxLevel);
279    % determine the mutation level
280    mutationLevel = floor(log2(mutationPoint))+1;
281    % delete the original heap nodes
282    for i = 1:length(mutateLocations)
283    heap{mutateLocations(i)} = [];
284    end
285    if mutationLevel == maxLevel
286        heap(mutateLocations) = varConst(randi(length(varConst)));
287    else
288        % generate a subtree
289        subHeap = heapGeneration(operator, varConst, (maxLevel-mutationLevel+1));
290        for i = 1:length(mutateLocations)
291            heap(mutateLocations(i)) = subHeap(i);
292        end
293    end
294    % replace c with constant
295    heap = replaceC(heap);
296    end
297
298    function operateLocations = searchChildren(operatPoint, maxLevel)
299    searchQueue = [operatPoint];
300    operateLocations = [operatPoint];
301    while ~isempty(searchQueue)
302        currentIndex = searchQueue(1);
303        if currentIndex*2+1 <= 2^maxLevel - 1
304            operateLocations(end+1) = currentIndex*2;
305            operateLocations(end+1) = currentIndex*2+1;
306            searchQueue(end+1) = currentIndex*2;
307            searchQueue(end+1) = currentIndex*2+1;
308        end
309        searchQueue(1) = [];
310    end
311    end
312
313    function [] = randomSearch(operator, varConst, maxLevel, data, n, repeat)
314    parfor r = 1:repeat
315        fileID = fopen(strcat('Random_',int2str(r),'.txt'),'wt');
316        bestError = inf;
317        %bestHeap = cell([2^maxLevel-1,1]);
318        for i = 1:n
319            % maximum level of heap can vary from 2 to 4
320            heap = heapGeneration(operator, varConst, maxLevel);
321            error = MAE_Cal(heap, data);
322            if error < bestError
323                bestError = error;
324                %bestHeap = heap;
325            end
326            fprintf(fileID, '%d %10.8f %10.8f \n', i, error, bestError);
327        end
328        fprintf("RM: Eval: %d, Error: %10.8f.\n", n, bestError)
329    end
330    end
331
332    function eval = evalHeap(heap, data)
333    heapStr = heapString(heap);
334    fh = str2func(heapStr);
335    eval = fh(data(:,1));
```

14

```
336  end
337
338  function error = MAE_Cal(heap, data)
339  heapStr = heapString(heap);
340  fh = str2func(heapStr);
341  y = fh(data(:,1));
342  error = sum(abs(data(:,2) - y))/length(y);
343  end
344
345  function heap = heapGeneration(operator, varConst, maxLevel)
346  heapSize = 2^maxLevel - 1;
347  heap = cell([heapSize,1]);
348  opSize = 2^(maxLevel-2)-1;
349  operatorQueue = [1];
350  while ~isempty(operatorQueue)
351      % pick the current index
352      currentIndex = operatorQueue(1);
353      % assign current operator
354      heap(currentIndex) = operator(randi(length(operator)));
355      % make sure the operator assignment does not exceed the limit
356      if currentIndex <= opSize
357          if rand < 0.5
358              operatorQueue(end+1) = currentIndex*2;
359          end
360          if rand < 0.5
361              operatorQueue(end+1) = currentIndex*2+1;
362          end
363      end
364      % delete the first in queue
365      operatorQueue(1) = [];
366  end
367  % record the operator index
368  opIndex = find(~cellfun(@isempty,heap));
369  varConstIndex = [];
370  for i = 1:length(opIndex)
371      if isempty(heap{opIndex(i)*2})
372          varConstIndex(end+1) = opIndex(i)*2;
373      end
374      if isempty(heap{opIndex(i)*2+1})
375          varConstIndex(end+1) = opIndex(i)*2+1;
376      end
377  end
378  % assign x and c to the rest of the tree
379  for i = 1:length(varConstIndex)
380      heap(varConstIndex(i)) = varConst(randi(length(varConst)));
381  end
382  % replace the 'c' with actual constant
383  heap = replaceC(heap);
384  end
385
386  function heap = replaceC(heap)
387  heapSize = length(heap);
388  for i = 1:heapSize
389      const = -10:0.1:10;
390      if strcmp(heap{i},'c')
391          heap{i} = const(randi(length(const)));
392      end
393  end
394  end
395
396  function heapStr = heapString(heap)
397  heapStr = heap;
398  % reverse order from 15 to 1
399  for i = fliplr(1:floor(length(heap)/2))
400      % make sure current node is not empty
401      if ~isempty(heapStr{i})
402          % if current node is sin or cos, only combine its right child
403          if strcmp(heapStr{i},'sin') || strcmp(heapStr{i},'cos')
```

15

```matlab
            %heapStr{i} = strcat(num2str(heapStr{i*2}),'*',num2str(heapStr{i}),'(',num2str(
                heapStr{2*i+1}),')');
            %heapStr{i} = strcat(num2str(heapStr{i}),'(',num2str(heapStr{2*i+1}),')');
            heapStr{i} = strcat('(',num2str(heapStr{i}),'(',num2str(heapStr{2*i}),'.*',num2str(
                heapStr{2*i+1}),')',')');
        else
            % else combine the left child, current node and right child
            heapStr{i} = strcat('(',num2str(heapStr{2*i}),num2str(heapStr{i}),num2str(heapStr{2*i
                +1}),')');
        end
    end
end
% take the top of the heap as output
heapStr = strcat('@(x) ',heapStr{1});
end
```