Saemundur Haraldsson John Woodward Sandy Brownlee



Computational Heuristics Operational Research Decision-Support

Fixing bugs in Python programs with Genetic Improvement

Program size and search granularity





Overview of talk

- Developing a GI framework for Python programs
- Search granularity and program size
- Breaking and fixing small Python programs

Motivation

- GI has already been successfully applied to large software, >50K LOC (Langdon et al. & Le Goues et al.)
- Pushing GI to its lower size limit for usefulness
- "The competent programmer hypothesis" for students
- Easier to analyse exactly what the GI is doing

GI for Python

GI for Python ----- Entities of the population

- Evolving Edit lists
 - A single edit: < "Edit", "Old code", "New code", "Location">
- Available edits
 - Copy, Swap, Delete and Replace
- Movable code
 - Whole Lines
 - **Boolean** operators: 'or', 'and', 'not', '<=', '!=', etc.
 - Mathematical operators: '+', '*', '-', '%', etc
 - **Incremental** operators: '+=', '*=', '/=', '-='
 - Numerical constants
- Fitness function
 - Number of passed test cases

- The usual customizable properties
 - Population size
 - Number of generations
 - Selection
 - Survival / Elitism
- Offspring entities made with mutation only
 - Grow: Append randomly generated edits
 - Prune: Shorten the list of edits
 - Single edit mutation: Randomly select 1 edit and change it slightly.



- The usual customizable properties
 - Population size
 - Number of generations
 - \circ Selection
 - Survival / Elitism
- Offspring entities made with mutation only
 - Grow: Append randomly generated edits
 - Prune: Shorten the list of edits
 - Single edit mutation: Randomly select 1 edit and change it slightly.



- The usual customizable properties
 - Population size
 - Number of generations
 - \circ Selection
 - Survival / Elitism
- Offspring entities made with mutation only
 - Grow: Append randomly generated edits
 - Prune: Shorten the list of edits
 - Single edit mutation: Randomly select 1 edit and change it slightly.



- The usual customizable properties
 - Population size
 - Number of generations
 - \circ Selection
 - Survival / Elitism
- Offspring entities made with mutation only
 - Grow: Append randomly generated edits
 - Prune: Shorten the list of edits
 - Single edit mutation: Randomly select 1 edit and change it



Search Granularity Program Size



Search Granularity ----- Experimental setup

Movable code

- Random line edits
- Like for like line edits
- Change operators: math, boolean and incremental.

Step size (mutation choices)

- Grow and prune only (variable size)
- Single edit mutations and Grow (single edit growth)
- Both above

		Step size					
	Movable code	All available	Grow and Prune	Single edit			
	Random lines			X			
	Like for like lines	X	X	X			
	Operators and numbers	X	X	X			

Program size

- Lines of Code
 - Ranging from 5 100
- Implemented from various online sources
 - "100+ python challenging programming exercises"
 - <u>www.ActiveState.com</u> -- code recipies
 - www.Cprogramming.com -- challenge
- Beginner level programs that contain common code elements
 - Simple numerical calculations: Factorial
 - Mathematical constants approximations: pi, e, sqrt(2)
 - Simple text input Calculator
 - etc.

Breaking and Fixing

Breaking and fixing, The breaking process



Breaking and fixing, The fixing process

- Objectives are:
 - Number of test cases passed
 - Size of edit list, i.e. number of changes to the broken program
- Runs for 50 generations (population of 20)
- Returns the overall best solution.
 - Fewest number of changes made to the program to pass the greatest number of test cases.



Experiments, Line for line

		Broken	Fixed		100 experiments
Program	Size LOC	Avg. size of breaker	Avg. evals -> fixed	Avg. proportion of error variants	Avg. size of fixer
count_digs_letters	9	1	15.2	<mark>75%</mark>	2.01
dict_square	5	1	6.3	<mark>68%</mark>	1.5
divisable_5	7	1	10.2	<mark>81%</mark>	3.7
even_digits	13	1	4	<mark>74%</mark>	1.2
factorial	<mark>5</mark>	N/A	N/A	<mark>100%</mark>	N/A
formula_this	8	1	6.2	<mark>72%</mark>	4.1

Experiments, Line for line

		Broken	Fixed		
Program	Size LOC	Avg. size of breaker	Avg. evals -> fixed	Avg. proportion of error variants	Avg. size of fixer
lines_2_list	12	1	10.9	<mark>67%</mark>	4.01
list_tuple	5	<mark>N/A</mark>	N/A	<mark>100%</mark>	N/A
make_multiMatrix	8	1	14.5	<mark>80%</mark>	3.4
sort_unique	5	1	13.2	<mark>45%</mark>	2.13
sort_words	5	1	8.4	<mark>51%</mark>	1.25

Experiments, Summary of line for line

- Breaking
 - Fitness is effectively binary: broken or not broken
 - pass all or no test cases
 - Highly unlikely programming errors.
 - e.g. forgetting a complete line?
 - Takes only one line out of place to break.
 - If a valid break exists it is found in first generation.
- Fixing
 - Takes longer to find the fix than the break
 - High proportion of variants do not run
 - and those that run are mostly semantically identical, i.e. loads of redundancy

Experiments, finer grained

Case example, Dictionary of squares

- Input: single integer n
- Output: dictionary of all the numbers squared from 0 to n
- 5 test cases which include **boundary inputs**, n = 0 and 1
- Program was broken by replacing the first occurrence of 1 with 2.
 - o <REPLACE, '1', '2', 2,15>
- Then the GI was run 100 times to fix.
 - No elitism





Experiments, Finer grained: Dictionary of squares



Experiments, finer grained: Dictionary of squares

Experiments, finer grained

Case example: A simple text input calculator

- ~100 LOC
- Inserted bugs with 4 edits
 - Forced by increasing the required failed test cases
- Fails all test cases (19)
 - At least one test case for each function: +, -, *, and /
 - and the rest combines them
- Again: GI run 100 times to fix
 - Now with elitism

Experiments, finer grained



Experiments, summary of finer grained

- Sometimes finds mutations that pass some test cases
 - Fitness is not always binary, rather a step: passes 1 or 2 boundary cases.
 - More bugs -> more needles
- Much more realistic programming errors
 - typing "=" instead of "+=" or "<" instead of "<="
- Only one edit needed to break



Experiments, summary of finer grained

- We can nearly always find a valid break
 - Syntactically correct programs
 - High proportion of variants run
- For such small programs the fix is usually converting it back to the original.
 - No clever fixes, that weren't foreseen.
- The fix is most often found in the first 5-10 generations.
- Still, finding the fix takes much longer than finding the break.
 - In practice "Needle/s in a haystack" fitness function that is largely level.



Summary

- GI for Python programs is doable and promising
- Tested on multiple small programs
- Considered 2 dimensions of search granularity
 - Step size
 - Movable code
- Line based GI is not a realistic option for small programs
 - Where the boundary of size lies remains to be confirmed
- Smaller programs call for finer grained searches



Computational Heuristics Operational Research Decision-Support

Thanks for listening Questions?





