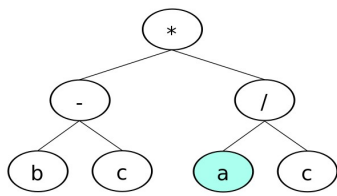


# Karoo GP User Guide

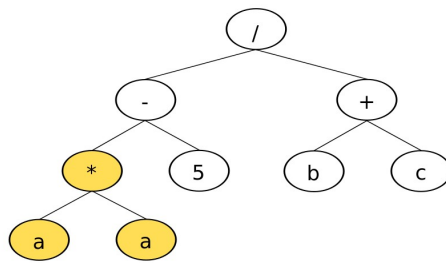
## Genetic Programming in Python

ver. 20160724

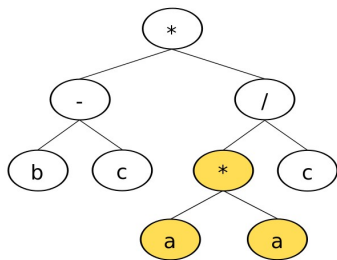
by Kai Staats



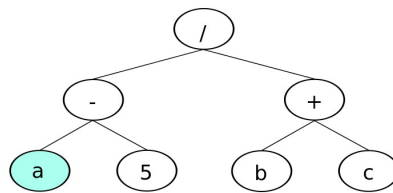
Parent A



Parent B



Child A



Child B

```

Tree 14 yields (sym): a - a/c - b*c**2 + b*c + b - 2*c - 1/c + c**2/b + c/b + 1/b
Tree 15 yields (sym): -a**3 - a*b**2 + a/(b*c**2) - 4*b + 3*c + 1 - 1/c
Tree 16 yields (sym): a**3*c + a*b + a*b/c - a + b*c + 3*c + 1 - 1/c + c/b
Tree 17 yields (sym): -a**2*c/b - 2*a + a/b**2 - b**2 - 2*b*c - 3*b + c/b - c/a
Tree 18 yields (sym): a*b**2 + a*c - a - a/c - b**2 - b + 2*b/c - c - c/b - b**2/a + c**2/a
Tree 19 yields (sym): -2*a**2 + a*b - b + 2*c - 2 + 1/b - c/b**2 + c/a + 1/a
Tree 20 yields (sym): -a**2*c - a**2/c - a - a/b - b**3*c + b**2 + 2*b + c**2 + c - 1 - c**3/a
Tree 21 yields (sym): -a**2*c - a*b + a - a*c/b - a*c/b**2 - c**2 + c
Tree 22 yields (sym): a**2 - a*b*c - a*b + a*c + b*c**3 - 2*b/c - c**2 + 4*c - 1 + b*c/a
Tree 23 yields (sym): a*b*c + a*b + a*b/c - a + b**2*c + b*c**3 + b - c + 1
Tree 24 yields (sym): -a*b*c + a*c**2 + a*c - 2*a - b*c + b + 4*c - 1 + b/(a*c) + 1/(a*c)
Tree 25 yields (sym): -a*b*c**2 + a*b*c - a*c + 2*a + 3*b - 3*c - 1 + 1/b - b**2/(a*c)
Tree 26 yields (sym): -a*b*c + 2*a*b - b**2*c - b - c**2 - 3*c + 2 - c**2/b - b*c/a
Tree 27 yields (sym): -a**2 - a/c + a*c**2/b + a/b - a*c/b**3 + b*c - b + 2*c

```

## Table of Contents

Welcome	3
What is Genetic Programming	3
The Karoo GP Applications	4
First time through	5
Second time through	5
Subsequent runs	6
Kernel Select	6
Type of Tree	6
Base Tree Depth	7
Maximum Tree Depth	7
Minimum Number of Nodes	7
Number of Trees	7
Number of Generations	7
Display Mode	7
What you see on-screen	8
The Evolutionary Operators	9
Reproduction	9
Point Mutation	9
Branch Mutation	9
Crossover	10
Functions (Operators)	10
Terminals (Operands)	12
Data Management	13
Default Data Format	13
Loading Your Dataset	13
Data Size	13
Train vs Test	13
Tree Population Management	14
Feature Generation	15
Tools	16
Karoo Data Sort	16
Karoo Data Normalise	16
Karoo Multiclass Classifier	16
Notes	17

```
Tree 13 yields (sym): a - a/c - b*c**2 + b*c + b - 2*c - 1/c + c**2/b + c/b + 1/b
Tree 14 yields (sym): a - a/c - b*c**2 + b*c + b - 2*c - 1/c + c**2/b + c/b + 1/b
Tree 15 yields (sym): -a**3 - a*b**2 + a/(b*c**2) - 4*b + 3*c + 1 - 1/c
Tree 16 yields (sym): a**3*c + a*b + a*b/c - a + b*c + 3*c + 1 - 1/c + c/b
Tree 17 yields (sym): -a**2*c/b - 2*a + a/b**2 - b**2 - 2*b*c - 3*b + c/b - c/a
Tree 18 yields (sym): a*b**2 + a*c - a - a/c - b**2 - b + 2*b/c - c - c/b - b**2/a + c**2/a
Tree 19 yields (sym): -2*a**2 + a*b - b + 2*c - 2 + 1/b - c/b**2 + c/a + 1/a
Tree 20 yields (sym): -a**2*c - a**2/c - a - a/b - b**3*c + b**2 + 2*b + c**2 + c - 1 - c**3/a
Tree 21 yields (sym): -a**2*c - a*b + a - a*c/b - a*c/b**2 - c**2 + c
Tree 22 yields (sym): a**2 - a*b*c - a*b + a*c + b*c**3 - 2*b/c - c**2 + 4*c - 1 + b*c/a
Tree 23 yields (sym): a*b*c + a*b + a*b/c - a + b**2*c + b*c**3 + b - c + 1
Tree 24 yields (sym): -a*b*c + a*c**2 + a*c - 2*a - b*c + b + 4*c - 1 + b/(a*c) + 1/(a*c)
Tree 25 yields (sym): -a*b*c**2 + a*b*c - a*c + 2*a + 3*b - 3*c - 1 + 1/b - b**2/(a*c)
Tree 26 yields (sym): -a*b*c + 2*a*b - b**2*c - b - c**2 - 3*c + 2 - c**2/b - b*c/a
Tree 27 yields (sym): -a**2 - a/c + a*c**2/b + a/b - a*c/b**3 + b*c - b + 2*c
```

## Karoo GP User Guide

*Welcome to Karoo GP!*

Karoo GP is an evolutionary algorithm, a *genetic programming* application suite which provides both symbolic regression and classification analysis. Written in the programming language [Python](#), Karoo GP owes its foundation to the “[Field Guide to Genetic Programming](#)” by Poli, Langdon, McPhee, and Koza.

Karoo GP provides a transparent interface to the inner workings of tree-based Genetic Programming. As a teaching tool, it enables instructors to share step-by-step how an evolutionary algorithm arrives to its solution. As a hands-on learning tool, Karoo GP supports rapid, repeatable experimentation with a simple, no-programming-required interface. Included with Karoo GP are two executables: the Desktop with an intuitive text-based user interface; Server with a combination of a configuration script and support for command-line execution and flagging. Both work with the built-in, test cases or external, user provided data.

This guide is designed to give you a crash course in tree-based Genetic Programming and a light introduction to preparing datasets for machine learning. Even if you are an expert in machine learning, you are encouraged to make time to read this document in order to anticipate all that Karoo GP has to offer.

*Requirements:*

- Python 2.7.6 [GCC 4.8.2]
- libraries: [NumPy](#) 1.9.2, [SymPy](#) 0.7.6.1, [pprocess](#) 0.5.1, [sklearn](#) 0.16.1, [matplotlib](#), csv, os, sys, argparse, and time

*Tested systems:*

- Ubuntu Linux 14.04 Desktop on Intel i5/i7 (4 core): no known issues
- Ubuntu Linux 14.04.3 LTS on Intel Xeon E5-2650 (40 core): no known issues

## What is Genetic Programming?

Genetic Programming (GP) is a type of evolutionary algorithm, a subset of machine learning. Inspired by biological evolution, GP is a computer program that uses random mutation, crossover, a fitness function, and multiple generations of evolution to resolve a user-defined task. At the foundation level, genetic programming discovers the relationship between dataset features—correlations between measurements of the real world.

Long before the advent of gene sequencing or machine learning, taxonomists conducted precise measurements of physical plants and animals in order to differentiate species. In 1936, biologist R.A. Fisher worked to differentiate species of Iris flowers using data generated by Edgar Anderson. In his paper “The use of multiple measurements in taxonomic problems”, Fisher states “When two or more populations have been measured in several characters,  $x_1, \dots, x_8$ , special interest attaches to certain linear functions of the measurements by which the populations are best discriminated.” Using Genetic Programming, we readily discover the mathematical functions that cleanly separate the 3 Iris species using just 50 data points each. This Iris multivariate dataset is a simple yet classic, must-solve problem for all Machine Learning developers of classification algorithms, and is included with Karoo GP as the default Classification dataset.

In general, genetic programming can be used to discover a relationship between features in data (symbolic regression) or to group data into categories (classification). Both of these employ a training (learning) and test (validation) phase, after which the resulting function (mathematical expression) may be employed in a live data stream for realtime regression or classification analysis, accordingly.

```
Tree 13 yields (sym): a - b - 2*c - 2*b - 2*c/b - 2*c/a
Tree 14 yields (sym): a - a/c - b*c**2 + b*c + b - 2*c - 1/c + c**2/b + c/b + 1/b
Tree 15 yields (sym): -a**3 - a*b**2 + a/(b*c**2) - 4*b + 3*c + 1 - 1/c
Tree 16 yields (sym): a**3*c + a*b + a*b/c - a + b*c + 3*c + 1 - 1/c + c/b
Tree 17 yields (sym): -a**2*c/b - 2*a + a/b**2 - b**2 - 2*b*c - 3*b + c/b - c/a
Tree 18 yields (sym): a*b**2 + a*c - a - a/c - b**2 - b + 2*b/c - c - c/b - b**2/a + c**2/a
Tree 19 yields (sym): -2*a**2 + a*b - b + 2*c - 2 + 1/b - c/b**2 + c/a + 1/a
Tree 20 yields (sym): -a**2*c - a**2/c - a - a/b - b**3*c + b**2 + 2*b + c**2 + c - 1 - c**3/a
Tree 21 yields (sym): -a**2*c - a*b + a - a*c/b - a*c/b**2 - c**2 + c
Tree 22 yields (sym): a**2 - a*b*c - a*b + a*c + b*c**3 - 2*b/c - c**2 + 4*c - 1 + b*c/a
Tree 23 yields (sym): a*b*c + a*b + a*b/c - a + b**2*c + b*c**3 + b - c + 1
Tree 24 yields (sym): -a*b*c + a*c**2 + a*c - 2*a - b*c + b + 4*c - 1 + b/(a*c) + 1/(a*c)
Tree 25 yields (sym): -a*b*c**2 + a*b*c - a*c + 2*a + 3*b - 3*c - 1 + 1/b - b**2/(a*c)
Tree 26 yields (sym): -a*b*c + 2*a*b - b**2*c - b - c**2 - 3*c + 2 - c**2/b - b*c/a
Tree 27 yields (sym): -a**2 - a/c + a*c**2/b + a/b - a*c/b**3 + b*c - b + 2*c
```

## The Karoo GP Applications

### Desktop

The Karoo GP Desktop application presents a functional interface for configuring, launching, and monitoring your Karoo GP runs. There are 5 choices for how you monitor the evolutionary process, from the *silent* mode in which you see only the outcome of each generation to a detailed, step-by-step debug mode in which every mutation and crossover process is presented to you, as both a validation and learning tool.

To launch the Desktop application from the shell:

```
$ python karoo_gp_main.py
```

If you include the path to an external dataset, it will load at launch:

```
$ python karoo_gp_main.py /[path]/[to_your]/[filename].csv
```

While this can be run on a remote server, you may find that once you get the hang of using Karoo, and are in production mode, `karoo_gp_server.py` provides both a powerful, scripted and command-line platform.

### Server

Karoo GP Server is designed to support repeated, scripted, and command-line executed genetic programming runs. You can manually configure a routine to move through a series of tree depths or population sizes, or apply an external genetic algorithm to analyze the leading expressions from each generation and modify the configuration of each subsequent generation or full run accordingly, to help guide GP toward the optimal solution.

```
$ python karoo_gp_server.py
```

Without arguments, Karoo GP defaults to the datasets located in `karoo_gp/files/` as determined by the kernel selected, eg: a Classification kernel loads `data_CLASSIFY.csv`. If you include *only* a path to your own dataset:

```
$ python karoo_gp_server.py /[path]/[to_your]/[filename].csv
```

If, however, you include any other arguments, then you must *also* include a flag to load your own dataset:

```
$ python karoo_gp_main.py -ker c -typ r -bas 4 -fil /[path]/[to_your]/[filename].csv
```

You can set your preferred default configuration parameters in the file `karoo_gp_server.py` such that without any command-line arguments, Karoo GP Server will repeat a run for comparison to priors. Command-line arguments then override these default settings, as follows:

argument	options	description
-ker	[r,c,m]	kernel: (r)egression, (c)lassification, or (m)atching
-typ	[f,g,r]	tree type: (f)ull, (g)row, or (r)amped half/half
-bas	[3...10]	maximum tree depth for the initial population
-max	[3...10]	maximum tree depth for the entire run
-min	[3...100]	minimum number of nodes
-pop	[10...1000]	maximum population
-gen	[1...100]	number of generations
-fil	/[path]/[to_your]/[filename].csv	load an external dataset

# Karoo GP

```
Tree 13 yields (sym): a - b - 2*c - 2*b/a - 2*c/a
Tree 14 yields (sym): a - a/c - b*c**2 + b*c + b - 2*c - 1/c + c**2/b + c/b + 1/b
Tree 15 yields (sym): -a**3 - a*b**2 + a/(b*c**2) - 4*b + 3*c + 1 - 1/c
Tree 16 yields (sym): a**3*c + a*b + a*b/c - a + b*c + 3*c + 1 - 1/c + c/b
Tree 17 yields (sym): -a**2*c/b - 2*a + a/b**2 - b**2 - 2*b*c - 3*b + c/b - c/a
Tree 18 yields (sym): a*b**2 + a*c - a - a/c - b**2 - b + 2*b/c - c - c/b - b**2/a + c**2/a
Tree 19 yields (sym): -2*a**2 + a*b - b + 2*c - 2 + 1/b - c/b**2 + c/a + 1/a
Tree 20 yields (sym): -a**2*c - a**2/c - a - a/b - b**3*c + b**2 + 2*b + c**2 + c - 1 - c**3/a
Tree 21 yields (sym): -a**2*c - a*b + a - a*c/b - a*c/b**2 - c**2 + c
Tree 22 yields (sym): a**2 - a*b*c - a*b + a*c + b*c**3 - 2*b/c - c**2 + 4*c - 1 + b*c/a
Tree 23 yields (sym): a*b*c + a*b + a*b/c - a + b**2*c + b*c**3 + b - c + 1
Tree 24 yields (sym): -a*b*c + a*c**2 + a*c - 2*a - b*c + b + 4*c - 1 + b/(a*c) + 1/(a*c)
Tree 25 yields (sym): -a*b*c**2 + a*b*c - a*c + 2*a + 3*b - 3*c - 1 + 1/b - b**2/(a*c)
Tree 26 yields (sym): -a*b*c + 2*a*b - b**2*c - b - c**2 - 3*c + 2 - c**2/b - b*c/a
Tree 27 yields (sym): -a**2 - a/c + a*c**2/b + a/b - a*c/b**3 + b*c - b + 2*c
```

## First time through

The Karoo GP user interface will prompt you for a series of inputs. Your first time through, enter Play mode in order to understand the functions of Karoo GP and to learn to visualise a GP tree, as follows:

1. Select Play mode.
2. Select a Full tree.
3. Select a depth of 1.

Depth 1: 1 operator, 2 variables, eg:  $a / c$

Depth 2: 3 operators, 4 variables, eg:  $a + b - 2 * c$

Depth 3: 7 operators, 8 variables, eg:  $c + b * c^3 / a + c / a$

Tree Depth: 0 of 1

```
NODE: 1
type: root
label: /
arity: 2
parent node:
child node(s): 2 3
```

Tree Depth: 1 of 1

```
NODE: 2
type: term
label: a
arity: 0
parent node: 1
child node(s):
```

```
NODE: 3
type: term
label: c
arity: 0
parent node: 1
child node(s):
```

Tree 1 yields (raw):  $(a)/(c)$   
Tree 1 yields (sym):  $a/c$

Take a look the GP tree represented on-screen. There are 3 nodes composed of 2 terminals (variables) and 1 function (operator). Select Play mode a few more times, playing with larger trees and both Full and Grow methods. You will find that the branches of Full trees always reach the bottom, where the branches of Grow trees may terminate early, with a terminal always the final node. With some practice, you can read the tree as printed on screen. Future versions of Karoo GP will include an option to generate an image (similar to the cover page).

## Second time through

Run Karoo GP again, but this time simply press ENTER at each of the queries. Sit back and watch as Karoo GP works to resolve a simple relationship between 3 variables which represent 3 columns of data:  $a + b + c = s$ .

We know the answer, but Karoo GP must discover a solution through an evolutionary process. As GP is based on random number generation, there is a chance that in the default 10 generations it will not work, or it might find a relationship between the columns of data other than that which was expected.

By configuring a run for deeper trees, or by increasing the minimum number of nodes required, you can encourage the evolutionary process to discover more complex solutions.

When done, you will be presented with *pause*. Type ? and ENTER to review the menu:

1. Type *l* and ENTER to view all trees which provide the correct solution.
2. Type *p*, ENTER, and then the unique number of any tree in the list, to print that tree to screen.
3. The solution will be listed in both its *raw* and Sympyfied format. Remember that multiple trees may carry the same or more than one solution to the same problem.
4. You may adjust the (b)alance of the operators, (cont)inue the evolutionary process, (test) the trees, or (q)uit.

```
Tree 93 yields (sym): a + 2*b + 1
Tree 94 yields (sym): a + 2*b + c/b
Tree 95 yields (sym): 2*b + c/b**2
Tree 96 yields (sym): a + a/c + b + c
Tree 97 yields (sym): a + b + c + 1
Tree 98 yields (sym): a + a/c + b + c
Tree 99 yields (sym): a + b + b/c
Tree 100 yields (sym): a + 3*b
```

28 trees [ 1 2 3 4 5 6 7 8 9 10 14 18 21 25 37 38 39 49 50 55 59 64 74 76 84 85 91 93] offer the highest fitness scores.

```

Tree 13 yields (sym): a - a/c - b*c**2 + b*c + b - 2*c - 1/c + c**2/b + c/b + 1/b
Tree 14 yields (sym): a - a/c - b*c**2 + b*c + b - 2*c - 1/c + c**2/b + c/b + 1/b
Tree 15 yields (sym): -a**3 - a*b**2 + a/(b*c**2) - 4*b + 3*c + 1 - 1/c
Tree 16 yields (sym): a**3*c + a*b + a*b/c - a + b*c + 3*c + 1 - 1/c + c/b
Tree 17 yields (sym): -a**2*c/b - 2*a + a/b**2 - b**2 - 2*b*c - 3*b + c/b - c/a
Tree 18 yields (sym): a*b**2 + a*c - a - a/c - b**2 - b + 2*b/c - c - c/b - b**2/a + c**2/a
Tree 19 yields (sym): -2*a**2 + a*b - b + 2*c - 2 + 1/b - c/b**2 + c/a + 1/a
Tree 20 yields (sym): -a**2*c - a**2/c - a - a/b - b**3*c + b**2 + 2*b + c**2 + c - 1 - c**3/a
Tree 21 yields (sym): -a**2*c - a*b + a - a*c/b - a*c/b**2 - c**2 + c
Tree 22 yields (sym): a**2 - a*b*c - a*b + a*c + b*c**3 - 2*b/c - c**2 + 4*c - 1 + b*c/a
Tree 23 yields (sym): a*b*c + a*b + a*b/c - a + b**2*c + b*c**3 + b - c + 1
Tree 24 yields (sym): -a*b*c + a*c**2 + a*c - 2*a - b*c + b + 4*c - 1 + b/(a*c) + 1/(a*c)
Tree 25 yields (sym): -a*b*c**2 + a*b*c - a*c + 2*a + 3*b - 3*c - 1 + 1/b - b**2/(a*c)
Tree 26 yields (sym): -a*b*c + 2*a*b - b**2*c - b - c**2 - 3*c + 2 - c**2/b - b*c/a
Tree 27 yields (sym): -a**2 - a/c + a*c**2/b + a/b - a*c/b**3 + b*c - b + 2*c

```

## Subsequent runs

When you are ready to dive in a bit further, take a few minutes to learn about each of the following:

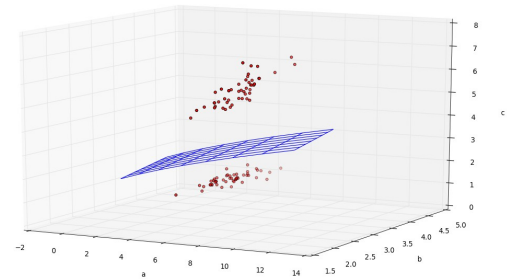
### 1 – Kernel Select

When prompted with a choice of (c)lassify, (r)egression, (m)atching, or (p)lay:

**Regression** is a minimisation function, meaning it seeks a fitness score with the *lowest* number. By default, this kernel will attempt to solve Kepler's 3rd Law of Planetary Motion. As the correct solution is  $t^2 / r^3$ , the challenge is for GP to not fall to a local minimum of  $t/t$  or  $r/r$ . You can experiment with both deeper trees and increasing the minimum number of nodes. Or, let Karoo GP run for 20, 30, ... 50 or more generations and random mutation might just discover the correct solution.

To learn more, visit: [www.physicsclassroom.com/class/circles/Lesson-4/Kepler-s-Three-Laws](http://www.physicsclassroom.com/class/circles/Lesson-4/Kepler-s-Three-Laws)

**Classification** is a maximisation function, meaning it seeks a fitness score with the *highest* number. By default, this kernel will attempt to solve the Iris flower problem. This is a machine learning classic built upon data generated before the dawn of DNA classification, when botanists used microscopes and calipers to measure the features of plants and animals. While Kepler's Law has one solution, the Iris problem has many. If you plot one of the GP generated functions over a scatter plot of the data (see tools/) you should find a clean separation of 2 or 3 species.



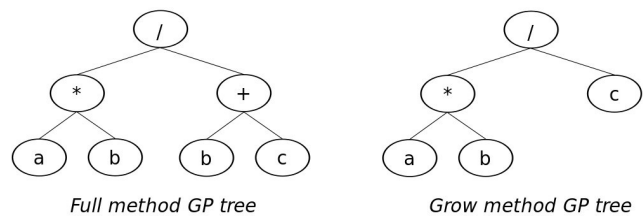
To learn more, visit: [archive.ics.uci.edu/ml/datasets/Iris](http://archive.ics.uci.edu/ml/datasets/Iris)

Karoo GP supports multiclass classification from 2 to  $n$  classes, with automated scaling. See *Additional Scripts* (below) to learn more, and to play with the included `karoo_multiclassifier.py` found in the `tools/` directory.

**Match** will attempt to discover a relationship between variables which produces an exact match, as with  $a + b + c = s$ , as noted in *Second Time Through* (above). This is the simplest of the fitness functions to experiment with on your own, but has limited realworld application. Try increasing tree depth and the minimum number of nodes. See *Using Your Own Data* at the bottom of this guide to learn how to replace the default data with your own.

### 2 – Type of Tree

All *Full* tree branches reach the maximum depth. These trees are less likely to solve problems as they are not as flexible in their solution space. For example, if the desired solution is  $a + b + c$  (5 elements) but a Full tree is confined to 15 elements, some will be forced to cancel each other in order to arrive to the desired solution. However, Full trees do contribute higher order expressions.



*Grow* trees have unbalanced branches, meaning some branches reach the maximum depth while others do not. *Grow* trees are more likely to find simpler solutions, but may not discover improved solutions in a higher order space.

Originated by John Koza, *Ramped Half/Half* initialises the first population with a 50/50 split of Full/Grow methods, and a spread of depths from min to max. The *Grow* method is then applied for each subsequent population. This is the most popular method as it injects a higher degree of diversity into the initial population.

```

Tree 13 yields (sym): a - a/c - b*c**2 + b*c + b - 2*c - 1/c + c**2/b + c/b + 1/b
Tree 14 yields (sym): a - a/c - b*c**2 + b*c + b - 2*c - 1/c + c**2/b + c/b + 1/b
Tree 15 yields (sym): -a**3 - a*b**2 + a/(b*c**2) - 4*b + 3*c + 1 - 1/c
Tree 16 yields (sym): a**3*c + a*b + a*b/c - a + b*c + 3*c + 1 - 1/c + c/b
Tree 17 yields (sym): -a**2*c/b - 2*a + a/b**2 - b**2 - 2*b*c - 3*b + c/b - c/a
Tree 18 yields (sym): a*b**2 + a*c - a - a/c - b**2 - b + 2*b/c - c - c/b - b**2/a + c**2/a
Tree 19 yields (sym): -2*a**2 + a*b - b + 2*c - 2 + 1/b - c/b**2 + c/a + 1/a
Tree 20 yields (sym): -a**2*c - a**2/c - a - a/b - b**3*c + b**2 + 2*b + c**2 + c - 1 - c**3/a
Tree 21 yields (sym): -a**2*c - a*b + a - a*c/b - a*c/b**2 - c**2 + c
Tree 22 yields (sym): a**2 - a*b*c - a*b + a*c + b*c**3 - 2*b/c - c**2 + 4*c - 1 + b*c/a
Tree 23 yields (sym): a*b*c + a*b + a*b/c - a + b**2*c + b*c**3 + b - c + 1
Tree 24 yields (sym): -a*b*c + a*c**2 + a*c - 2*a - b*c + b + 4*c - 1 + b/(a*c) + 1/(a*c)
Tree 25 yields (sym): -a*b*c**2 + a*b*c - a*c + 2*a + 3*b - 3*c - 1 + 1/b - b**2/(a*c)
Tree 26 yields (sym): -a*b*c + 2*a*b - b**2*c - b - c**2 - 3*c + 2 - c**2/b - b*c/a
Tree 27 yields (sym): -a**2 - a/c + a*c**2/b + a/b - a*c/b**3 + b*c - b + 2*c

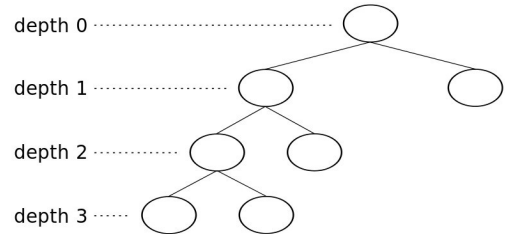
```

### 3 – Base Tree Depth

This is the tree depth for the *initial* population. By default, this is also the Maximum Tree Depth (below), unless otherwise user specified. If the base and maximum tree depths are equal, this will inhibit bloat, but also restrict potential, more complex solutions.

The depth of a tree is illustrated to the right, and relates to the maximum possible number of nodes, as follows:

$$\text{nodes} = 2^{(\text{depth} + 1)} - 1$$



### 4 – Maximum Tree Depth

Karoo is unique from traditional tree-based GP in that it incorporates a user defined maximum tree depth, thereby restricting program bloat. However, deeper trees present opportunity for more complex solutions, and they enable the inclusion of a greater number of features (columns in your .csv). As depth 3 enables up to 15 nodes, and depth 5 enables up to 63 nodes, quite a bit of growth is possible with trees just 2-3 depths greater. Or, you can set the maximum to 10 (2047 possible nodes) and learn if a larger solution is in fact more able to resolve your data.

### 5 – Minimum Number of Nodes

If the Maximum Tree Depth is the *ceiling*, then this is the *floor*. The minimum number of nodes (both operators and terminals) defines the simplest expression allowed. For example, the correct solution to Kepler's 3<sup>rd</sup> Law of Planetary Motion is  $t * t / r * r * r$  which has 9 elements (nodes). In Karoo GP, the gene pool from which the tournament selection operates is built only from those trees which meet the minimum number of nodes criterion. This is very useful in solving the default problem for the Regression kernel, as GP tends to the simpler  $t / t$ .

**NOTE:** If you set the minimum number of nodes too high, you may invoke elitism in the population, killing off simpler solutions or even the entire population.

### 6 – Number of Trees

100 trees per population is a good place to start. However, larger populations offer a greater diversity of possible solutions. Feel free to experiment with larger populations.

### 7 – Number of Generations

For simpler problems with less than a dozen features, 10 generations will often give you a sense of whether or not Karoo GP is on the right track. With more complex problems, 20, 30, ... 50 or more generations may be required for the random process of evolution to generate an improved tree which provides the desired solution.

You may (cont)inue with subsequent generations when all generations have run their course (see bottom).

### 8 – Display Mode

The display modes (i)nteractive, (m)inimal, (g)eneration, and (s)ilent are the means by which you interact with and monitor the evolutionary process of Karoo GP. They do not affect the outcome of the process itself.

**(i)nteractive** - *pause* with each step of the Karoo GP process.

**(m)inimal** - runs through the entire evolution, start to finish without *pause*. However, it does display each tree as it is evaluated, which is quite entertaining (if you are into that sort of thing).

**(g)eneration** - *pause* with the end of each generation, allowing you to make adjustments once per generation.

# Karoo GP

```
Tree 13 yields (sym): a - a/c - b*c**2 + b*c + b - 2*c - 1/c + c**2/b + c/b + 1/b
Tree 14 yields (sym): a - a/c - b*c**2 + b*c + b - 2*c - 1/c + c**2/b + c/b + 1/b
Tree 15 yields (sym): -a**3 - a*b**2 + a/(b*c**2) - 4*b + 3*c + 1 - 1/c
Tree 16 yields (sym): a**3*c + a*b + a*b/c - a + b*c + 3*c + 1 - 1/c + c/b
Tree 17 yields (sym): -a**2*c/b - 2*a + a/b**2 - b**2 - 2*b*c - 3*b + c/b - c/a
Tree 18 yields (sym): a*b**2 + a*c - a - a/c - b**2 - b + 2*b/c - c - c/b - b**2/a + c**2/a
Tree 19 yields (sym): -2*a**2 + a*b - b + 2*c - 2 + 1/b - c/b**2 + c/a + 1/a
Tree 20 yields (sym): -a**2*c - a**2/c - a - a/b - b**3*c + b**2 + 2*b + c**2 + c - 1 - c**3/a
Tree 21 yields (sym): -a**2*c - a*b + a - a*c/b - a*c/b**2 - c**2 + c
Tree 22 yields (sym): a**2 - a*b*c - a*b + a*c + b*c**3 - 2*b/c - c**2 + 4*c - 1 + b*c/a
Tree 23 yields (sym): a*b*c + a*b + a*b/c - a + b**2*c + b*c**3 + b - c + 1
Tree 24 yields (sym): -a*b*c + a*c**2 + a*c - 2*a - b*c + b + 4*c - 1 + b/(a*c) + 1/(a*c)
Tree 25 yields (sym): -a*b*c**2 + a*b*c - a*c + 2*a + 3*b - 3*c - 1 + 1/b - b**2/(a*c)
Tree 26 yields (sym): -a*b*c + 2*a*b - b**2*c - b - c**2 - 3*c + 2 - c**2/b - b*c/a
Tree 27 yields (sym): -a**2 - a/c + a*c**2/b + a/b - a*c/b**3 + b*c - b + 2*c
```

**(s)ilent** - designed to run remotely on a server with limited overhead. This mode presents the results of each generation just to let you know it is still alive 'n cranking through data while you are engaged in data reductions, Facebook, or a Star Wars role playing game.

There are 2 additional modes available to you at any time, as listed with the *pause* menu:

**(d)e(b)ug** - for those of you who want to watch the evolutionary process at the intimate level, this is a fully transparent window to the inner workings of tree-by-tree Mutation and Crossover. In addition, you will watch the criteria of the engaged fitness function, the tournament selection, and the population of the gene pool from which the next generation will be selected.

**(t)imer** - allows you to monitor the performance of your current GP run against the number of CPU cores currently engaged. At each *pause*, you may select (c), enter the number of CPU cores, and then ENTER again to continue. Karoo GP is *multi-core* engaged, meaning it spawns processes across physical CPU cores (not multiple threads on a single core). At the time of this edit, Karoo GP is *not* GPU enabled.

You can switch to any mode at any *pause*. However, once you are in *minimal* or *silent* modes, there are no longer opportunities to modify the evolutionary parameters until all generations are complete. However, you may continue the evolutionary process and Karoo GP will pick up where you left off. Simply enter (cont) followed by the number of subsequent generations, and off you go. If new parameters are desired, remember to set all configurations before you continue.

## What You See On-Screen

### The *pause* menu

Be certain to play around with all the options available to you in *pause*. "?" or "help" followed by ENTER will display all available options. You can change the level of interaction, minimum number of nodes, number of CPU cores; print and test trees, and more.

### The most fit trees

At the end of each generation of population evolution, a list of bold numbers is presented on-screen. These are the trees GP has found to offer the best overall fitness scores in the latest generation. At each *pause*, you may (l)ist, (p)rint, and/or *test* (various tests may be applied) any given tree, and again at the end of any given run.

The leading Trees and their associated expressions are:

1 :  $2*a - b + 2*c$

2 :  $a + 2*b + 1$

5 :  $a + 2*b + 1$

6 :  $a + b + c$

10 :  $a + 2*b + 1$

12 :  $a + 2*b + 1$

83 :  $2*a - b + 2*c$

21 :  $2*a - b + 2*c$

94 :  $a + 2*b + 1$



```

Tree 14 yields (sym): a - a/c - b*c**2 + b*c + b - 2*c - 1/c + c**2/b + c/b + 1/b
Tree 15 yields (sym): -a**3 - a*b**2 + a/(b*c**2) - 4*b + 3*c + 1 - 1/c
Tree 16 yields (sym): a**3*c + a*b + a*b/c - a + b*c + 3*c + 1 - 1/c + c/b
Tree 17 yields (sym): -a**2*c/b - 2*a + a/b**2 - b**2 - 2*b*c - 3*b + c/b - c/a
Tree 18 yields (sym): a*b**2 + a*c - a - a/c - b**2 - b + 2*b/c - c - c/b - b**2/a + c**2/a
Tree 19 yields (sym): -2*a**2 + a*b - b + 2*c - 2 + 1/b - c/b**2 + c/a + 1/a
Tree 20 yields (sym): -a**2*c - a**2/c - a - a/b - b**3*c + b**2 + 2*b + c**2 + c - 1 - c**3/a
Tree 21 yields (sym): -a**2*c - a*b + a - a*c/b - a*c/b**2 - c**2 + c
Tree 22 yields (sym): a**2 - a*b*c - a*b + a*c + b*c**3 - 2*b/c - c**2 + 4*c - 1 + b*c/a
Tree 23 yields (sym): a*b*c + a*b + a*b/c - a + b**2*c + b*c**3 + b - c + 1
Tree 24 yields (sym): -a*b*c + a*c**2 + a*c - 2*a - b*c + b + 4*c - 1 + b/(a*c) + 1/(a*c)
Tree 25 yields (sym): -a*b*c**2 + a*b*c - a*c + 2*a + 3*b - 3*c - 1 + 1/b - b**2/(a*c)
Tree 26 yields (sym): -a*b*c + 2*a*b - b**2*c - b - c**2 - 3*c + 2 - c**2/b - b*c/a
Tree 27 yields (sym): -a**2 - a/c + a*c**2/b + a/b - a*c/b**3 + b*c - b + 2*c
    
```

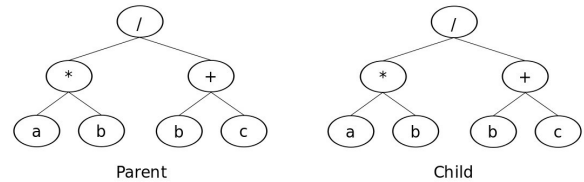
## The Evolutionary Operators

There are 4 evolutionary operators applied in Karoo GP. The foundation for these evolutionary operators was derived from the “[Field Guide to Genetic Programming](#)”\* by Riccardo Poli, William Langdon, and Nicholas McPhee, with contributions by John Koza. Using the (b)alance option in the *pause* menu, you can set the ratio of the operators, such that combined they equal 100%.

In the following, *tournament selection* refers to the random selection of a number of trees (default 10) whose fitness scores are compared. The tree with the highest score is moved to the next generation through an evolutionary operation. We use tournament selection, not a top-to-bottom evaluation of the entire population, else we risk *elitism*, premature convergence on what may not be the best overall solution.

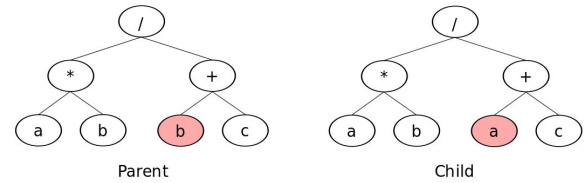
### Reproduction

Through tournament selection, a single tree from the prior population is copied without mutation to the next generation. In the biological world, this is analogous to a member of a population entering the gene pool of the subsequent (younger) generation.



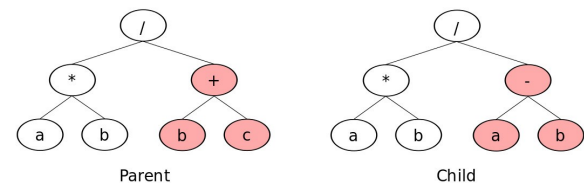
### Point Mutation

Through tournament selection, a copy of a tree from the prior population mutates a single node before being added to the next generation. In the biological world, this may be analogous to asexual reproduction, that is, a copy of an individual with a minor mutation. In this method, a single point is selected for mutation while maintaining function nodes as operators and terminal nodes as terminals. The size and shape of the tree will remain identical.



### Branch Mutation

Through tournament selection, a copy of a tree from the prior population mutates before being added to the next generation. In the biological world, this may be analogous to asexual reproduction, that is, a copy of an individual but with a potentially *substantial* mutation. Unlike Point Mutation, in this method an entire branch is selected. If the evolutionary run is designated as Full, the size and shape of the tree will remain identical, each node mutated sequentially, where functions remain functions and terminals remain terminals. If the evolutionary run is designated as Grow or Ramped Half/Half, the size and shape of the tree may grow smaller or larger, but it may not exceed the maximum depth defined by the user.



\* The Field Guide and many other books about GP are listed at [www.geneticprogramming.com](http://www.geneticprogramming.com)

# Karoo GP

```

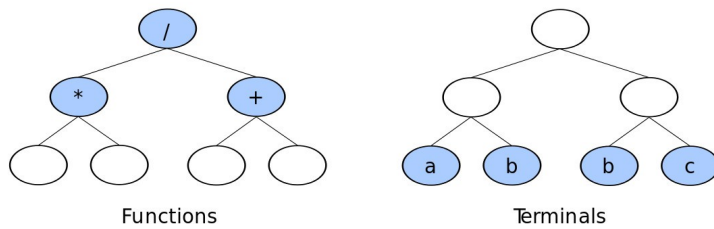
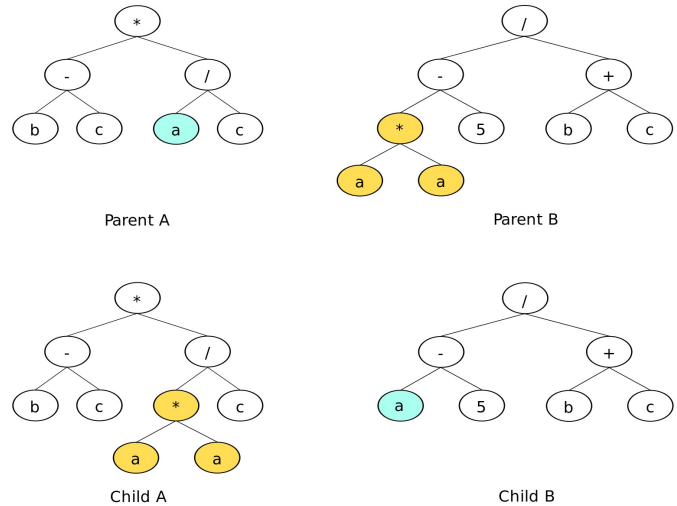
Tree 14 yields (sym): a - a/c - b*c**2 + b*c + b - 2*c - 1/c + c**2/b + c/b + 1/b
Tree 15 yields (sym): -a**3 - a*b**2 + a/(b*c**2) - 4*b + 3*c + 1 - 1/c
Tree 16 yields (sym): a**3*c + a*b + a*b/c - a + b*c + 3*c + 1 - 1/c + c/b
Tree 17 yields (sym): -a**2*c/b - 2*a + a/b**2 - b**2 - 2*b*c - 3*b + c/b - c/a
Tree 18 yields (sym): a*b**2 + a*c - a - a/c - b**2 - b + 2*b/c - c - c/b - b**2/a + c**2/a
Tree 19 yields (sym): -2*a**2 + a*b - b + 2*c - 2 + 1/b - c/b**2 + c/a + 1/a
Tree 20 yields (sym): -a**2*c - a**2/c - a - a/b - b**3*c + b**2 + 2*b + c**2 + c - 1 - c**3/a
Tree 21 yields (sym): -a**2*c - a*b + a - a*c/b - a*c/b**2 - c**2 + c
Tree 22 yields (sym): a**2 - a*b*c - a*b + a*c + b*c**3 - 2*b/c - c**2 + 4*c - 1 + b*c/a
Tree 23 yields (sym): a*b*c + a*b + a*b/c - a + b**2*c + b*c**3 + b - c + 1
Tree 24 yields (sym): -a*b*c + a*c**2 + a*c - 2*a - b*c + b + 4*c - 1 + b/(a*c) + 1/(a*c)
Tree 25 yields (sym): -a*b*c**2 + a*b*c - a*c + 2*a + 3*b - 3*c - 1 + 1/b - b**2/(a*c)
Tree 26 yields (sym): -a*b*c + 2*a*b - b**2*c - b - c**2 - 3*c + 2 - c**2/b - b*c/a
Tree 27 yields (sym): -a**2 - a/c + a*c**2/b + a/b - a*c/b**3 + b*c - b + 2*c
    
```

## Crossover

Through tournament selection, 2 trees are selected as *parents* to produce 2 *offspring*. Within each parent tree a branch is selected. For child A, parent A is copied, with its selected branch deleted. Parent B's branch is then copied to the former location of parent A's branch, and inserted (grafted). This is reversed for child B. The size and shape of the offspring may be smaller or larger than either of the parents, but may not exceed the maximum depth defined by the user.

This process combines genetic code from trees which were chosen by the tournament process as having a higher fitness than the average population. Therefore, there is a higher probability their offspring will provide an improved fitness.

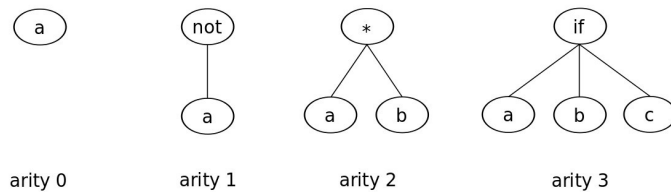
According to the literature, crossover is the most commonly applied evolutionary operator at 70-90%.



## Functions (Operators)

The operators applied to Karoo GP are designated by you. The name of the .csv file which contains your desired operators must match the name of the kernel chosen. For example, the Classification kernel will auto-load /files/functions\_CLASSIFY.csv.

The most commonly used are arithmetic [+ , - , \* , /], trigonometric [cos, sin], and Boolean [and, or]. In addition, [\*\*, sqrt, exp] are supported. Many more have not been tested, but are open to exploration. As all Karoo GP expressions are processed by the library Sympy, you can learn about all possible operators at [www.sympy.org](http://www.sympy.org)



All operators presented to Karoo GP via the associated must be followed by an *arity*, that is, the number of variables or constants each operator expects to work with. For example, \* has an arity of 2, meaning it expects something before and after itself, as in the expression  $a * b$

# Karoo GP

```
Tree 14 yields (sym): a - a/c - b**c**2 + b*c + b - 2*c - 1/c + c**2/b + c/b + 1/b
Tree 15 yields (sym): -a**3 - a*b**2 + a/(b*c**2) - 4*b + 3*c + 1 - 1/c
Tree 16 yields (sym): a**3*c + a*b + a*b/c - a + b*c + 3*c + 1 - 1/c + c/b
Tree 17 yields (sym): -a**2*c/b - 2*a + a/b**2 - b**2 - 2*b*c - 3*b + c/b - c/a
Tree 18 yields (sym): a*b**2 + a*c - a - a/c - b**2 - b + 2*b/c - c - c/b - b**2/a + c**2/a
Tree 19 yields (sym): -2*a**2 + a*b - b + 2*c - 2 + 1/b - c/b**2 + c/a + 1/a
Tree 20 yields (sym): -a**2*c - a**2/c - a - a/b - b**3*c + b**2 + 2*b + c**2 + c - 1 - c**3/a
Tree 21 yields (sym): -a**2*c - a*b + a - a*c/b - a*c/b**2 - c**2 + c
Tree 22 yields (sym): a**2 - a*b*c - a*b + a*c + b*c**3 - 2*b/c - c**2 + 4*c - 1 + b*c/a
Tree 23 yields (sym): a*b*c + a*b + a*b/c - a + b**2*c + b*c**3 + b - c + 1
Tree 24 yields (sym): -a*b*c + a*c**2 + a*c - 2*a - b*c + b + 4*c - 1 + b/(a*c) + 1/(a*c)
Tree 25 yields (sym): -a*b*c**2 + a*b*c - a*c + 2*a + 3*b - 3*c - 1 + 1/b - b**2/(a*c)
Tree 26 yields (sym): -a*b*c + 2*a*b - b**2*c - b - c**2 - 3*c + 2 - c**2/b - b*c/a
Tree 27 yields (sym): -a**2 - a/c + a*c**2/b + a/b - a*c/b**3 + b*c - b + 2*c
```

The power operator (\*\*) is not typically necessary as this is generated by any variable multiplied by itself. In the real world, most non-linear laws of nature are based on powers of 2 or 3, not ^5 or ^12. As such, GP finds these without the explicit operator. Square root too can be discovered through \* and /, but this may be introduced if anticipated.

As presented in karoo\_gp/files/operators\_list.txt, [cos, sin, exp, sqrt, log] must be preceded by another operator, as shown below. So, you will likely want to include all 4 of the arithmetic operators to provide a balanced set.

```
+ cos,2
- cos,2
* cos,2
/ cos,2
```

**NOTE:** *There is no space between the operator, the comma, and the arity.*

An example of the contents of a potential functions\_[kernel].csv file follows:

```
operator, arity
+,2
-,2
*,2
/,2
+,2
-,2
*,2
/,2
+,2
-,2
*,2
/,2
+,2
-,2
*,2
/,2
+,2
-,2
*,2
/,2
+ cos,2
- cos,2
* cos,2
/ cos,2
```

**NOTE:** *The selection of any operator from this list is entirely random. The quantity of times any given operator appears relates to its chance of being selected, and therefore biases the outcome of the evolutionary process. If you have a single trigonometric operator which is represented 4 times, then you may desire to provide each of the arithmetic operators 4 times as well.*

```

Tree 14 yields (sym): a - a/c - b*c**2 + b*c + b - 2*c - 1/c + c**2/b + c/b + 1/b
Tree 15 yields (sym): -a**3 - a*b**2 + a/(b*c**2) - 4*b + 3*c + 1 - 1/c
Tree 16 yields (sym): a**3*c + a*b + a*b/c - a + b*c + 3*c + 1 - 1/c + c/b
Tree 17 yields (sym): -a**2*c/b - 2*a + a/b**2 - b**2 - 2*b*c - 3*b + c/b - c/a
Tree 18 yields (sym): a*b**2 + a*c - a - a/c - b**2 - b + 2*b/c - c - c/b - b**2/a + c**2/a
Tree 19 yields (sym): -2*a**2 + a*b - b + 2*c - 2 + 1/b - c/b**2 + c/a + 1/a
Tree 20 yields (sym): -a**2*c - a**2/c - a - a/b - b**3*c + b**2 + 2*b + c**2 + c - 1 - c**3/a
Tree 21 yields (sym): -a**2*c - a*b + a - a*c/b - a*c/b**2 - c**2 + c
Tree 22 yields (sym): a**2 - a*b*c - a*b + a*c + b*c**3 - 2*b/c - c**2 + 4*c - 1 + b*c/a
Tree 23 yields (sym): a*b*c + a*b + a*b/c - a + b**2*c + b*c**3 + b - c + 1
Tree 24 yields (sym): -a*b*c + a*c**2 + a*c - 2*a - b*c + b + 4*c - 1 + b/(a*c) + 1/(a*c)
Tree 25 yields (sym): -a*b*c**2 + a*b*c - a*c + 2*a + 3*b - 3*c - 1 + 1/b - b**2/(a*c)
Tree 26 yields (sym): -a*b*c + 2*a*b - b**2*c - b - c**2 - 3*c + 2 - c**2/b - b*c/a
Tree 27 yields (sym): -a**2 - a/c + a*c**2/b + a/b - a*c/b**3 + b*c - b + 2*c

```

## Terminals (Operands)

Karoo GP automatically extracts the terminals (variables and constants) from the top of the data .csv file which you have provided. This enables you to name each column of your data to represent that particular feature. Later, when you read a GP derived expression, you will see how it relates to the features and original data.

## Variables

If you are attempting to predict the weather using historical data<sup>[1]</sup> you might have something like this:

jd	sat	saz	hum	tmp	wm	wd	s
0	0.9292	1	0.2613	0.549	0.3956	0.0837	0
0.0092	0.9027	0.9865	0.2368	0.5862	0.4377	0.1328	0
0.0184	0.8761	0.9697	0.195	0.5908	0.4048	0.0913	0
0.0276	0.8496	0.9562	0.2021	0.6077	0.1905	0.1159	0
0.0368	0.8142	0.9461	0.2076	0.6156	0.2106	0.1221	0
0.0461	0.7788	0.9327	0.1821	0.611	0.185	0.0999	0
0.0553	0.7522	0.9226	0.1993	0.6099	0.152	0.3873	0
0.0645	0.7168	0.9125	0.2109	0.6043	0.2271	0.1416	0

Where “jd” is the variable for “Julian Day” (time), “sat” for the Sun's altitude, “saz” for the Sun's azimuth, “hum” for relative humidity, “tmp” for temperature, “wm” for the magnitude of the wind, and “wd” for the direction of the wind. The final column, “s” is the solution, that is, what you are solving for. In this case, we are applying a regression algorithm, with a minimisation function approaching zero in order to discover a correlation between the given variables.

## Constants

You can inject constants into the mix. In the current implementation of Karoo GP, constants are added at the top of unique columns in the data .csv. The entire column associated with each constant is then filled with zeros. These zeros are space fillers which have no impact on the processing of the data or operator.

For example, you might add a set of constant values [-5,-4,-3,-2,-1,1,2,3,4,5] which are balanced across the origin. Or, if you have a sense of how the variables need to be weighted, you may choose a particular set of constants. The same weather dataset above could include 5 constants expressed as decimal fractions, as the predicted variations in weights applied to this normalised data is relatively small.

jd	sat	saz	hum	tmp	wm	wd	0.1	0.2	0.3	0.4	0.5	s
0	0.9292	1	0.2613	0.549	0.3956	0.0837	0	0	0	0	0	0
0.0092	0.9027	0.9865	0.2368	0.5862	0.4377	0.1328	0	0	0	0	0	0
0.0184	0.8761	0.9697	0.195	0.5908	0.4048	0.0913	0	0	0	0	0	0
0.0276	0.8496	0.9562	0.2021	0.6077	0.1905	0.1159	0	0	0	0	0	0
0.0368	0.8142	0.9461	0.2076	0.6156	0.2106	0.1221	0	0	0	0	0	0
0.0461	0.7788	0.9327	0.1821	0.611	0.185	0.0999	0	0	0	0	0	0
0.0553	0.7522	0.9226	0.1993	0.6099	0.152	0.3873	0	0	0	0	0	0
0.0645	0.7168	0.9125	0.2109	0.6043	0.2271	0.1416	0	0	0	0	0	0

Future versions of Karoo GP will include constants as defined by the user via the interface.

[1] data from the SALT weather monitoring station and weather prediction software developed by Dr. Stephen Potter, head astronomer, South African Astronomical Observatory; Karoo GP is used to determine the weights

```
Tree 13 yields (sym): a - a/c - b*c**2 + b*c + b - 2*c - 1/c + c**2/b + c/b + 1/b
Tree 14 yields (sym): a - a/c - b*c**2 + b*c + b - 2*c - 1/c + c**2/b + c/b + 1/b
Tree 15 yields (sym): -a**3 - a*b**2 + a/(b*c**2) - 4*b + 3*c + 1 - 1/c
Tree 16 yields (sym): a**3*c + a*b + a*b/c - a + b*c + 3*c + 1 - 1/c + c/b
Tree 17 yields (sym): -a**2*c/b - 2*a + a/b**2 - b**2 - 2*b*c - 3*b + c/b - c/a
Tree 18 yields (sym): a*b**2 + a*c - a - a/c - b**2 - b + 2*b/c - c - c/b - b**2/a + c**2/a
Tree 19 yields (sym): -2*a**2 + a*b - b + 2*c - 2 + 1/b - c/b**2 + c/a + 1/a
Tree 20 yields (sym): -a**2*c - a**2/c - a - a/b - b**3*c + b**2 + 2*b + c**2 + c - 1 - c**3/a
Tree 21 yields (sym): -a**2*c - a*b + a - a*c/b - a*c/b**2 - c**2 + c
Tree 22 yields (sym): a**2 - a*b*c - a*b + a*c + b*c**3 - 2*b/c - c**2 + 4*c - 1 + b*c/a
Tree 23 yields (sym): a*b*c + a*b + a*b/c - a + b**2*c + b*c**3 + b - c + 1
Tree 24 yields (sym): -a*b*c + a*c**2 + a*c - 2*a - b*c + b + 4*c - 1 + b/(a*c) + 1/(a*c)
Tree 25 yields (sym): -a*b*c**2 + a*b*c - a*c + 2*a + 3*b - 3*c - 1 + 1/b - b**2/(a*c)
Tree 26 yields (sym): -a*b*c + 2*a*b - b**2*c - b - c**2 - 3*c + 2 - c**2/b - b*c/a
Tree 27 yields (sym): -a**2 - a/c + a*c**2/b + a/b - a*c/b**3 + b*c - b + 2*c
```

## Data Management

### Default Data Format

While Karoo GP can be recoded to work with any data format, by default, it anticipates .csv (comma separated values) files as readily exported from spreadsheet applications. Each row in a data .csv file is a data point (eg: a stamp in a time series), with 2 or more columns which hold raw data values or a values derived through *feature extraction* (more below).

In supervised machine learning, it is standard for the right-most column to be the solution, either a value for which a regression algorithm is working to approach, or a label for each of 2 or more classes. Karoo GP anticipates the header for this column to be labeled lower-case “s”.

With the Iris dataset (included with Karoo GP as a demonstration of multi-class classification), each row represents a single flower in the real world. The columns hold the measurements for each flower. The right-most column holds the solution, in this case, the classification [0,1,2] for each flower. No matter if you are working with datasets from biology, sociology, or astrophysics, the function of the rows and columns remain the same.

### Loading your Dataset

Karoo GP is designed for easily repeatable experimentation. When you launch Karoo GP, data and function files associated with the user selected kernel are auto-loaded. For example, if you select “c” for classification, then “files/data\_CLASSIFY.csv” and the list of operators (see section *Functions (Operators)*) “files/functions\_CLASSIFY.csv” are loaded accordingly.

When you are ready to apply your own dataset, you can replace the data in one of the .csv files in the files/ directory (keeping the file name the same) or point Karoo GP to a properly prepared .csv which contains your data, as described in the Karoo GP Applications sections for Desktop and Server, above.

**NOTE:** At the time of this writing, run-time loading a dataset does not load a unique, associated set of operators. Rather, you must first edit files/functions\_[kernel].csv before your Karoo GP run. This is best done from a shell editor such as nano or VI.

A few things to keep in mind:

1. Learn about the Sort and Normalise scripts (see Tools below) included with this package which will help you properly prepare your own dataset.
2. If you generate your own dataset using a spreadsheet application, be certain to save the original spreadsheet as .ods or .xls in addition to the .csv required by Karoo GP in order to preserve any formula, notes, or formatting for reference or later modification.
3. If you replace the contents of an existing .csv in the files directory with your own, pay attention to the following:
  - If editing the .csv using a text editor, do NOT add spaces between the variables and the comma.
  - Label each column with a variable. For example: *a,b,c* or *a1,a2,a3* or give each column a variable which represents the real-world measurement, such as *t* for time or *p* for pressure.
  - The right-most column must be labeled “s” (lower case, no spaces before or after).
4. Be very wary of your editor applying hidden characters in the .csv as these will cause Karoo to crash.
5. If sending a .csv file through email, wrap it in an archive (.zip, .tar.gz) else data corruption might occur.

**NOTE:** Your .csv must conform to the format described in the above or bad things will happen.

```

Tree 14 yields (sym): a - a/c - b*c**2 + b*c + b - 2*c - 1/c + c**2/b + c/b + 1/b
Tree 15 yields (sym): -a**3 - a*b**2 + a/(b*c**2) - 4*b + 3*c + 1 - 1/c
Tree 16 yields (sym): a**3*c + a*b + a*b/c - a + b*c + 3*c + 1 - 1/c + c/b
Tree 17 yields (sym): -a**2*c/b - 2*a + a/b**2 - b**2 - 2*b*c - 3*b + c/b - c/a
Tree 18 yields (sym): a*b**2 + a*c - a - a/c - b**2 - b + 2*b/c - c - c/b - b**2/a + c**2/a
Tree 19 yields (sym): -2*a**2 + a*b - b + 2*c - 2 + 1/b - c/b**2 + c/a + 1/a
Tree 20 yields (sym): -a**2*c - a**2/c - a - a/b - b**3*c + b**2 + 2*b + c**2 + c - 1 - c**3/a
Tree 21 yields (sym): -a**2*c - a*b + a - a*c/b - a*c/b**2 - c**2 + c
Tree 22 yields (sym): a**2 - a*b*c - a*b + a*c + b*c**3 - 2*b/c - c**2 + 4*c - 1 + b*c/a
Tree 23 yields (sym): a*b*c + a*b + a*b/c - a + b**2*c + b*c**3 + b - c + 1
Tree 24 yields (sym): -a*b*c + a*c**2 + a*c - 2*a - b*c + b + 4*c - 1 + b/(a*c) + 1/(a*c)
Tree 25 yields (sym): -a*b*c**2 + a*b*c - a*c + 2*a + 3*b - 3*c - 1 + 1/b - b**2/(a*c)
Tree 26 yields (sym): -a*b*c + 2*a*b - b**2*c - b - c**2 - 3*c + 2 - c**2/b - b*c/a
Tree 27 yields (sym): -a**2 - a/c + a*c**2/b + a/b - a*c/b**3 + b*c - b + 2*c

```

## Data Size

In theory, the size of the .csv file is limited only by the amount of RAM in your computer ... and your patience. A practical run of GP is often against 10,000 rows of data in a given .csv file. You can imagine that this takes some time, far more than the simple Kepler and Iris examples given.

## Train vs Test

Using a SciKit Learn module, Karoo GP auto-splits all datasets greater than 10 lines into TRAINING and TEST in order to offer data on which to test that is unique from the training. However, if you present Karoo with 10 or fewer rows in your .csv, it automatically copies the same data into both TRAIN and TEST as it is assumed you are using Karoo GP in an instructional or experimental mode. This will not provide a valid run.

## Tree Population Management

All trees generated for all populations of any given Karoo GP run are stored in files/ as follows:

*population\_a.csv* - includes the foundation population and all subsequent evolved populations. If you run 10 generations with 100 trees each, you will have a .csv file with 1000 trees.

*population\_b.csv* - written only if you export a snapshot of the evolutionary process (at the *pause* menu) using the (w)rite command.

*population\_final.csv* - the final generation of GP trees should, if all goes well, offer the most fit trees of all the generations. While this final generation is also included in *population\_a*, this file is isolated in order to make it easier for you to locate and work with these solutions.

TREE_ID	100								
tree_type	g								
tree_depth_base	4								
NODE_ID	1	2	3	4	5	6	7	8	9
node_depth	0	1	1	2	2	3	3	3	3
node_type	root	func	term	func	func	term	term	term	term
node_label	*	*	b	+	-	b	b	c	a
node_parent		1	1	2	2	4	4	5	5
node_arity	2	2	0	2	2	0	0	0	0
node_c1	2	4		6	8				
node_c2	3	5		7	9				
node_c3									
fitness	5			19					

## Archiving and Loading Seed Populations

To archive a population of trees at any time, copy files/population\_f.csv to a safe place.

When you desire to restore a population, rename the prior as population\_s.csv and move it back into the files/ directory. From the *pause* menu, select "load" and immediately, the former population\_s will be copied to *and overwrite* population\_a which is the foundation population for the next generation.

With Karoo GP you can load a former population or a hand-crafted, *seed population* at any *pause*. This enables you to work with archived populations to carry forward with an evolutionary run with a new set of parameters to test a new hypothesis, work against a new dataset, or conduct a new series of validations. What's more, you can start with the same population and test any number of evolutionary parameter configurations over varying generations in order to determine which enables the most rapid convergence on a desired solution.

If you desire to load an archived population select the (g)eneration or (i)nteractive display mode. Just after launch, at the first *pause*, conduct the load, configure your run, and then press ENTER to continue.

**NOTE:** The Karoo GP population files are very specific in their format. Do not alter the file structure, in any way.

```
Tree 13 yields (sym): a - a/c - b*c**2 + b*c + b - 2*c - 1/c + c**2/b + c/b + 1/b
Tree 14 yields (sym): a - a/c - b*c**2 + b*c + b - 2*c - 1/c + c**2/b + c/b + 1/b
Tree 15 yields (sym): -a**3 - a*b**2 + a/(b*c**2) - 4*b + 3*c + 1 - 1/c
Tree 16 yields (sym): a**3*c + a*b + a*b/c - a + b*c + 3*c + 1 - 1/c + c/b
Tree 17 yields (sym): -a**2*c/b - 2*a + a/b**2 - b**2 - 2*b*c - 3*b + c/b - c/a
Tree 18 yields (sym): a*b**2 + a*c - a - a/c - b**2 - b + 2*b/c - c - c/b - b**2/a + c**2/a
Tree 19 yields (sym): -2*a**2 + a*b - b + 2*c - 2 + 1/b - c/b**2 + c/a + 1/a
Tree 20 yields (sym): -a**2*c - a**2/c - a - a/b - b**3*c + b**2 + 2*b + c**2 + c - 1 - c**3/a
Tree 21 yields (sym): -a**2*c - a*b + a - a*c/b - a*c/b**2 - c**2 + c
Tree 22 yields (sym): a**2 - a*b*c - a*b + a*c + b*c**3 - 2*b/c - c**2 + 4*c - 1 + b*c/a
Tree 23 yields (sym): a*b*c + a*b + a*b/c - a + b**2*c + b*c**3 + b - c + 1
Tree 24 yields (sym): -a*b*c + a*c**2 + a*c - 2*a - b*c + b + 4*c - 1 + b/(a*c) + 1/(a*c)
Tree 25 yields (sym): -a*b*c**2 + a*b*c - a*c + 2*a + 3*b - 3*c - 1 + 1/b - b**2/(a*c)
Tree 26 yields (sym): -a*b*c + 2*a*b - b**2*c - b - c**2 - 3*c + 2 - c**2/b - b*c/a
Tree 27 yields (sym): -a**2 - a/c + a*c**2/b + a/b - a*c/b**3 + b*c - b + 2*c
```

## Feature Generation

If a dataset is a collection of raw data points, one or more values or measurements related to observation of the real world, then features are *transformed data*, values extracted from the observational data. However, in the world of Machine Learning, a set of features is often referred to as a “dataset” even though Machine Learning seldom works with raw, observational data.

## Feature Selection

*Feature selection* is the process of selecting a subset of relevant features which best describe the dataset as a whole, as some features may not be rich in discriminatory power. Feature selection helps reduce the dimensionality of the data in order that the machine learning algorithm is more readily able to learn.

## Feature Extraction

*Feature extraction* works with observational data to build new, derived, informative values from or relationships between selected data parameters.

According to Dr. Arun Aniyan, post-doctorate fellow at the Square Kilometer Array, features are:

- Informative and non-redundant.
- Support subsequent learning and generalization.
- Reduce dimensionality (in some cases).
- Lead to better human interpretation of the overall data.

SciKit Learn offers freely available methods for conducting [Feature Selection](#) and [Feature Extraction](#) in Python.

## Write Your Own Fitness Function

Fitness functions lie at the heart of all machine learning. They may be very simple, classifying according to a value greater than or equal to zero as 1, or less than zero as 0. Or they may be very complex, applying external scripts which track trends in the classification process itself, applying statistical analysis to determine how to adjust the fitness functions parameters over time.

Advanced users may modify existing fitness function kernels or create their own. This guide does not provide a detailed guide, rather a pointer to the places in Karoo GP which require modification, as follows:

- In `karoo_gp_main.py` and/or `karoo_gp_server.py`, add the new kernel to the query for `gp.kernel`
- In `karoo_gp_base_class.py`:
  1. In the method `fx_karoo_data_load`, modify the dictionaries: `data_dict`, `func_dict`, and `fitt_dict`
  2. In the method `fx_fitness_gym`, search for “[other]” and modify.
  3. In the method `fx_fitness_eval`, search for “[other]” and modify.
  4. Now build your new fitness function as a Python method. Search for “`def fx_fitness_function_[other]`” to find the place-holder. Keep in mind that due to the multi-core functionality of this particular section, it is imperative that your sum be passed back through the `pprocess` portal, not saved as a global variable. Refer to the other fitness functions for examples.

```
Tree 13 yields (sym): a - a/c - b*c**2 + b*c + b - 2*c - 1/c + c**2/b + c/b + 1/b
Tree 14 yields (sym): a - a/c - b*c**2 + b*c + b - 2*c - 1/c + c**2/b + c/b + 1/b
Tree 15 yields (sym): -a**3 - a*b**2 + a/(b*c**2) - 4*b + 3*c + 1 - 1/c
Tree 16 yields (sym): a**3*c + a*b + a*b/c - a + b*c + 3*c + 1 - 1/c + c/b
Tree 17 yields (sym): -a**2*c/b - 2*a + a/b**2 - b**2 - 2*b*c - 3*b + c/b - c/a
Tree 18 yields (sym): a*b**2 + a*c - a - a/c - b**2 - b + 2*b/c - c - c/b - b**2/a + c**2/a
Tree 19 yields (sym): -2*a**2 + a*b - b + 2*c - 2 + 1/b - c/b**2 + c/a + 1/a
Tree 20 yields (sym): -a**2*c - a**2/c - a - a/b - b**3*c + b**2 + 2*b + c**2 + c - 1 - c**3/a
Tree 21 yields (sym): -a**2*c - a*b + a - a*c/b - a*c/b**2 - c**2 + c
Tree 22 yields (sym): a**2 - a*b*c - a*b + a*c + b*c**3 - 2*b/c - c**2 + 4*c - 1 + b*c/a
Tree 23 yields (sym): a*b*c + a*b + a*b/c - a + b**2*c + b*c**3 + b - c + 1
Tree 24 yields (sym): -a*b*c + a*c**2 + a*c - 2*a - b*c + b + 4*c - 1 + b/(a*c) + 1/(a*c)
Tree 25 yields (sym): -a*b*c**2 + a*b*c - a*c + 2*a + 3*b - 3*c - 1 + 1/b - b**2/(a*c)
Tree 26 yields (sym): -a*b*c + 2*a*b - b**2*c - b - c**2 - 3*c + 2 - c**2/b - b*c/a
Tree 27 yields (sym): -a**2 - a/c + a*c**2/b + a/b - a*c/b**3 + b*c - b + 2*c
```

## Tools

The following scripts are found inside the included directory `karoo_gp/tools/`

### Karoo Data Sort

In machine learning, it is often the case that your engaged dataset is derived from a larger one. In constructing the subset, if we grab a series of data points (rows in a `.csv`) from the larger dataset in sequential order, only from the top, middle, or bottom, we will likely bias the new dataset and incorrectly train the machine learning algorithm. Therefore, it is imperative that we engage a random array, guided only by the number of data points for each class.

*This script can be used **\*before\*** `karoo_normalise.py`, and assumes no header has yet been applied to the `.csv`.*

### Karoo Data Normalise

This script works with a raw dataset to prepare a new, normalised dataset. It does so by comparing all values in each given column, finding the maximum and minimum values, and then modifying each value to fall between a high of 1 and low of 0. The modified values are written to a new file, the original remaining untouched.

*This script can be used **\*after\*** `karoo_features_sort.py`, and assumes no header has yet been applied to the `.csv`.*

### Karoo Multiclass Classifier

This is a toy script, designed to allow you to play with multiclass classification using the same underlying function as employed by Karoo GP. You can enter finite or infinite *wing* bins, and as many classes as you desire:

```
class 0 as -1.0 <= 0
class 0 as -0.5 <= 0
class 0 as 0.0 <= 0
class 1 as 0 < 0.5 <= 1
class 1 as 0 < 1.0 <= 1
class 2 as 1.5 > 1
class 2 as 2.0 > 1
class 2 as 2.5 > 1
```

### Karoo Iris Plot

This is a functional, even if rudimentary script designed to help you visualise your 2D or 3D data against an expression generated by Karoo GP. The challenge comes with solving complex equations for a single variable such that you have a plot-able expression. If the algebra required is beyond your skills (or you forgot what you learned in high school), tools such as Matlab may be of some assistance. If you desire to normalise your data in advance of using this script, the Karoo GP normalisation script included in the `karoo_gp/tools/` directory is very easy to use.

By default, this script plots a Karoo GP derived expression against a scatter plot of one of the Iris datasets included with this package: `karoo_gp/files/Iris_dataset/data_IRIS_virginica-vs-setosa_3-col_PLOT.csv`

## Experiment! Explore! *Evolve!*

You cannot damage Karoo GP by experimenting. If you screw up a data file, either Karoo GP or Python will complain (or just come to a screeching halt). Keep backups of your `data_` and `function_` files. Save templates of those that work well. But most of all, have fun! And if you find any errors in the code (there may be a few), do not hesitate to contact me. I will do my best to respond quickly –kai



# Karoo GP

```
Tree 13 yields (sym): a - a/c - b**2 + b - 2*c - 1/c + c**2/b + c/b + 1/b
Tree 14 yields (sym): a - a/c - b**c**2 + b*c + b - 2*c - 1/c + c**2/b + c/b + 1/b
Tree 15 yields (sym): -a**3 - a*b**2 + a/(b*c**2) - 4*b + 3*c + 1 - 1/c
Tree 16 yields (sym): a**3*c + a*b + a*b/c - a + b*c + 3*c + 1 - 1/c + c/b
Tree 17 yields (sym): -a**2*c/b - 2*a + a/b**2 - b**2 - 2*b*c - 3*b + c/b - c/a
Tree 18 yields (sym): a*b**2 + a*c - a - a/c - b**2 - b + 2*b/c - c - c/b - b**2/a + c**2/a
Tree 19 yields (sym): -2*a**2 + a*b - b + 2*c - 2 + 1/b - c/b**2 + c/a + 1/a
Tree 20 yields (sym): -a**2*c - a**2/c - a - a/b - b**3*c + b**2 + 2*b + c**2 + c - 1 - c**3/a
Tree 21 yields (sym): -a**2*c - a*b + a - a*c/b - a*c/b**2 - c**2 + c
Tree 22 yields (sym): a**2 - a*b*c - a*b + a*c + b*c**3 - 2*b/c - c**2 + 4*c - 1 + b*c/a
Tree 23 yields (sym): a*b*c + a*b + a*b/c - a + b**2*c + b*c**3 + b - c + 1
Tree 24 yields (sym): -a*b*c + a*c**2 + a*c - 2*a - b*c + b + 4*c - 1 + b/(a*c) + 1/(a*c)
Tree 25 yields (sym): -a*b*c**2 + a*b*c - a*c + 2*a + 3*b - 3*c - 1 + 1/b - b**2/(a*c)
Tree 26 yields (sym): -a*b*c + 2*a*b - b**2*c - b - c**2 - 3*c + 2 - c**2/b - b*c/a
Tree 27 yields (sym): -a**2 - a/c + a*c**2/b + a/b - a*c/b**3 + b*c - b + 2*c
```

## Notes