

2-bit Flip Mutation Elementary Fitness Landscapes

W. B. Langdon

Presented at Dagstuhl Seminar 10361, Theory of Evolutionary Algorithms, 8 September 2010

Fax: +44 (0)171 387 1397

Electronic Mail: W.Langdon@cs.ucl.ac.uk

URL: <http://www.cs.ucl.ac.uk/staff/W.Langdon/>

Abstract

Genetic Programming parity with only XOR is not elementary. GP parity can be represented as the sum of $k/2 + 1$ elementary landscapes. Statistics, including fitness distance correlation (FDC), of Parity's fitness landscape are calculated. Using Walsh analysis the eigen values and eigenvectors of the Laplacian of the two bit flip fitness landscape are given. Tests support λ/d as a measure of the ruggedness of elementary landscapes for predicting problem difficulty. An elementary needle in a haystack (NIH) landscape is given.

Keywords

genetic algorithms, genetic programming, search, optimisation, graph theory, Laplacian, Hamming cube

*Department of Computer Science
University College London
Gower Street
London WC1E 6BT, UK*

1 Fitness Landscapes

An elementary landscape is a special case of a fitness landscape. Firstly we recap fitness landscapes then some well known properties of elementary landscapes will be described in the next section. Section 3 describes our simplified version of the genetic programming parity problem. In Section 4 we reformulate tree GP parity as a bit string genetic algorithm where mutation flips exactly two bits. Section 5 gives many properties of the fitness landscapes created by two bit flip mutation, particularly those relating to GP parity. Section 6 describes some GP experiments on parity and on two elementary landscapes.

Fitness landscapes have often been used to try and explain how optimisation techniques, particularly evolutionary algorithms [1] such as genetic programming [11] work. An optimisation problem can be viewed as a search problem where all possible solutions to the problem are nodes in the search space and each has a value. In genetic algorithms [13] this is called a fitness value (more generally an objective value). Optimisation is viewed as sampling from this space with a goal to finding better points (or even the best point) in the space.

Except for Monte Carlo methods, optimisation techniques use information gathered from previous samples to decide where in the search space to sample next. The goal being to minimise the number of samples that are needed before an acceptable solution is found. Evolutionary algorithms, and several other optimisation techniques, only use the current search point (or the current population of search points) to guide the choice of where to look next. They do not use previously gained knowledge. Different algorithms can have radically different ways of moving from one point in the search space to the next. A search neighbourhood is the set of points that a specific algorithm can reach in one step from the current search point. A fitness landscape can be thought of as a graph where neighbours are linked by a single edge if and only if our search algorithm can move between the two nodes in the graph. While the height of the node is given by its fitness. See Figure 1.

Typically in evolutionary algorithms, the edges in the graph are undirected, because if one node can be reached from another then the reverse move is also possible. Notice that the difficulty of a problem depends not only on how fitness values are decided but also on the way the search algorithm moves across the search space. I.e. problem difficulty also depends on the fitness landscape the search algorithm imposes on the underlying problem [9, Chp. 2].

In genetic algorithms a common fitness landscape is to encode candidate solutions as strings of l bits. Each of the 2^l bit strings is allocated a fitness value. This gives a search space of 2^l candidate solutions. If mutation is restricted to flipping a single bit then each of the candidate solutions is connected in the fitness landscape to l other candidate solutions. This is known as the Hamming neighbourhood or hypercube graph [18]. The fitness landscape metaphor can be extended to population approaches (such as Particle Swarm Optimisation [10]) by allowing multiple sample points (one for each member of the population) in the landscape.

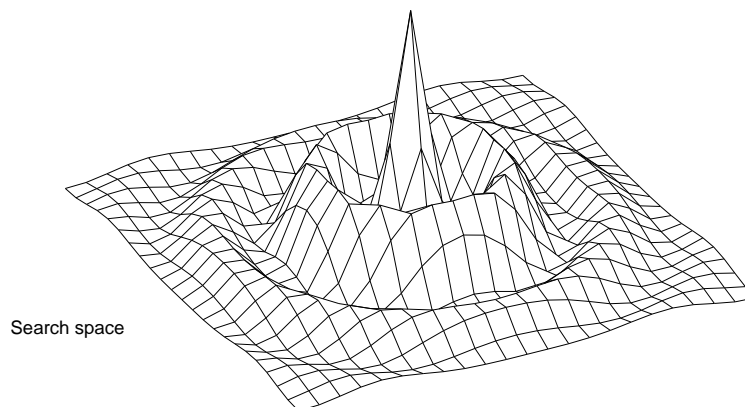


Figure 1: A fitness landscape where internal nodes have four neighbours. Fitness is plotted vertically.

Fitness landscapes are a useful metaphor with simple mutation which has a well defined neighbourhood. However with the canonical mutation only genetic algorithm the analogy starts to fail. Since GA mutation is defined as a probability of each bit flipping, multiple bits can be changed. There is a finite, albeit exponentially small, probability of any string being converted into any other. Thus, in principle, the whole fitness landscape becomes a single fully connected clique.

In principle fitness landscapes can be extended to allow multi-parent search operations (e.g. crossover) however then the neighbourhood of each member of the population depends upon the location of everyone else in the population and the idea of fixed predetermined links in the graph fails. (However it can be extended, e.g. via P-structures [16].)

2 Elementary Landscapes

2.1 Wave Equation

The following is based on Whitley's (e.g. [20]) definitions. (See also Stadler [15, p3].)

In an elementary landscape the fitness function f and the search space neighbourhood graph are related by a wave equation. The search space neighbourhood graph can be represented by a matrix Δ (see Sections 2.2 and 5.3). Only in an elementary landscape do the search space and the fitness function obey

$$\Delta f = \lambda(f - \bar{f})$$

Where \bar{f} is the mean fitness across the whole search space. That is, zero-mean fitness $f - \bar{f}$ is an eigenvector of Δ with eigenvalue λ . Notice a given search space (e.g. that created by one point mutation acting on a bit string chromosome) will have multiple eigenvectors, each of which will be the fitness function (up to an additive constant) for a different elementary landscape.

2.2 Average Fitness Change

In general in any regular landscape (i.e. any landscape where every location has the same number of search neighbours d) the mean change in fitness caused by one genetic change depends on the current position in the search space. Treating it as a vector gives:

$$\frac{1}{d}(Af - f)$$

If we treat the search space as a graph, then d is the degree (i.e. the number of links from the node) of each node in the graph. A is the adjacency matrix. For our purposes, it is a sparse real square matrix where every element corresponding to a link in the graph has the value 1 and all the others are zero.

$$\frac{1}{d}(Af - f) = \frac{1}{d}(A - dI)f = -\frac{1}{d}\Delta f$$

Where Δ is the Laplacian and is defined to be $dI - A$.

2.3 Average Neighbourhood Fitness in an Elementary Landscape

The mean fitness of a neighbourhood, $N(x)$, will also depend on the current position in the search landscape x . It is given by the fitness of the current position, $f(x)$, plus the average change in fitness (calculated in the previous section).

$$\begin{aligned}
 \text{avg}_{y \in N(x)} \{f(y)\} &= \frac{1}{d} \sum_{y \in N(x)} f(y) \\
 &= f(x) + \frac{1}{d} \sum_{y \in N(x)} f(y) - f(x) \\
 &= f(x) - \frac{1}{d} \Delta f(x) \\
 &= f(x) - \frac{1}{d} \lambda (f(x) - \bar{f}) \\
 &= f(x) + \frac{\lambda}{d} (\bar{f} - f(x))
 \end{aligned}$$

If, for convenience we set \bar{f} to zero, in an elementary landscape the mean fitness of the neighbours of x is:

$$\text{avg}_{y \in N(x)} \{f(y)\} = f(x) - \frac{\lambda}{d} f(x) = \left(1 - \frac{\lambda}{d}\right) f(x) \quad (1)$$

Thus, if $\lambda < d$, the mean of the neighbourhood is always closer to the overall average \bar{f} than the centre of the neighbourhood $f(x)$ is. Figure 2 shows a local optimum. By Equation 1, the average of the neighbourhood must lie between \bar{f} and $f(x)$. At a local optimum, by definition, $f(x)$ must be above the fitness of all its neighbours and hence must be above their average fitness. Therefore $f(x)$ must be above \bar{f} . Another way of saying this is: in elementary landscapes (with $\lambda < d$) there are no local optima with below average fitness.

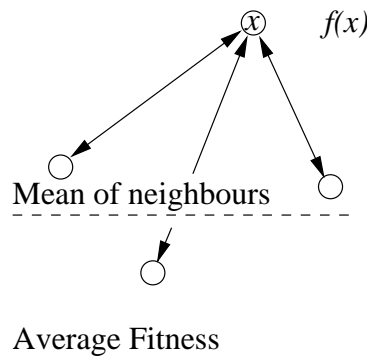


Figure 2: By definition at a hill top x the fitness of its neighbours is lower than that of the peak $f(x)$. By Equation 1, provided $\lambda < d$, the average fitness of the neighbours of x is closer to the average of the whole landscape \bar{f} than $f(x)$ is. Therefore $f(x)$ must lie above \bar{f} .

3 Genetic Programming Parity

The classic definition of the GP parity problems was given by Koza [6]. He defined the representation for order- k parity as binary trees whose external leafs are the functions inputs (drawn from $D_0 \cdots D_{k-1}$) and whose internal nodes are drawn from four binary Boolean functions. Initially the first simplification is to use just one Boolean function rather than four. We use either EQ or XOR depending if we are dealing with even or odd parity.

Koza defines the fitness of each tree by testing it on all 2^k possible input patterns and counting the number of times it returns the same answer as parity would. (Even- k -parity is true if the number of inputs which are true is even.) Thus Koza's fitness lies in the range $0 \dots 2^k$ and a solution to the parity problems has fitness 2^k .

Elsewhere [8, pages 421–423] we have used the symmetries of EQ and XOR to show that with a function set composed only of either EQ or only of XOR, trees will have fitness of either 2^{k-1} or 2^k . These properties also carry over from tree based GP to Cartesian GP and have also been exploited by Yu and Miller [22]. Further a tree is only a solution to parity if it contains an odd number of each type of leaf. (To simplify the text, and since odd parity behaves similarly, from now on we will discuss only even parity.) Since solutions to k -parity have one leaf of each of the k types plus redundant pairs of the same leaf type, they always have $k + 2n$ leafs (for $n = 0, 1, 2, \dots$). So their total size is $2k + 4n - 1$. Notice neither the shape of the tree nor the order of its leafs matters. All binary trees formed by re-arrangements of the leafs and internal nodes of the tree have identical fitness. The group properties of EQ also mean any pairs of leafs of the same type can be removed. E.g., using $=$ to represent EQ, $(D_1 = (D_2 = D_2))$ is identical to the input D_1 on its own. So if there are an even number any type of leaf (e.g. D_2) then it is as if the input was not connected. A tree missing one or more inputs scores exactly half the maximum score on parity.

We extend our previous analysis of the parity problem [8], which described its fitness distribution, to consider parity's fitness landscape.

4 The Parity Problem

4.1 The Mutation Operator

To form a landscape, in addition to the representation and its associated fitness function we need at least one search operator to establish which representations are adjacent to each other. Koza's [6] subtree crossover is a bit complicated to start with so we shall instead start with a mutation operator.

The mutation operator changes exactly one thing in the tree. Since the internal nodes are always EQ, the only possible change is to convert one leaf from one input (D_i) to another (D_j with $j \neq i$). This does not change the size of the tree. In principle more complicated tree size changing mutation operators or indeed crossover operators might also be considered.

4.2 Mutation's Impact on Fitness

The effect of mutation is either to make no difference to fitness (i.e. remain at 2^{k-1}) or to increase it from 2^{k-1} to 2^k or reduce it from 2^k to 2^{k-1} . A solution to parity has an odd number of all k types of leafs. Thus replacing any leaf with another of a different type will always mean it is no longer a solution. So it is impossible to mutate a tree of fitness 2^k into another tree of fitness 2^k .

Now the chance of each of these three things happening depends only on the number of each type of leaf. It does not depend upon the order of the leafs or their placement on the tree. Since these permutations make no difference to either fitness or to any of the changes in fitness we ignore them and initially replace trees with k integer counts d_i . Obviously there is a constraint that all of the tree's leafs must be present. So for a tree of size $2k + 4n - 1$, $\sum_{i=0}^{k-1} d_i$ must be $k + 2n$.

A tree's fitness is only 2^k if the number of leafs of each type (d_i) is odd for all k types. We can further simplify the representation by replacing the d_i s with e_i s which are 1 if their corresponding d_i is even. Similarly we can simplify Koza's fitness function so that 2^{k-1} is replaced by 0 and 2^k by 1. Under the new scheme fitness = 1 if and only if all e_i are 0 and fitness = 0 otherwise.

Mutation cannot change the total number of leafs, therefore having chosen an initial tree, $\sum d_i$ is fixed. This is a constraint on d_i and so we only need to specify $k - 1$ values for the d_i . (The remaining one can be inferred from the need to supply all the tree's leafs.) In fact, the oddness or evenness of $(k - 1) d_i$ is enough for the odd or evenness of the remaining one to be inferred. Another way of looking at this is to say: although for convenience we have k e_i , mutation can only reach half of their 2^k possible values. The unreachable half of the 2^k values correspond to trees of the wrong size for any of them to be parity solutions and so have zero fitness. We will assume the tree size is $2k + 4n - 1$, which allows one solution in a space of 2^{k-1} .

Mutating a leaf from type i to type j means decrementing d_i and incrementing d_j . Provided $d_i > 0$, this is equivalent to inverting e_i and e_j . I.e. mutation flips exactly two of the k bits. The next section will describe how we ensure $d_i > 0$ so that all mutation are always possible and further they are all equally likely.

4.3 Large Trees Simplify Mutation Analysis

Now it will greatly simplify things later if we assume the type of input we are about to mutate is chosen uniformly at random from the k possible types and that the type of the replacement leaf is chosen uniformly at random from the $k - 1$ remaining types. Obviously this requires the tree to have at least one leaf of each of the k types.

To further simplify the analysis we shall assume that the trees are much bigger than the minimum size ($2k - 1$). So big in fact, that we can assume that the tree has more than one leaf of each type. Further we assume the tree has sufficiently large number of leafs that undirected random drift from an initial random starting point in the search landscape will never cause the number of leafs of any of the k types to fall to one. I.e. random drift will not approach the edge of the k -dimensional simplex. If evolution lasts for G generations then drift will change the number of leafs of a given type by about \sqrt{G} . So assuming the tree is also bigger than a constant multiple of $(2k - 1)\sqrt{G}$ will ensure the chance of any of the k types of leaf approaching extinction is negligible. So mutation is always free to choose any pair of leaf types. This ensures it is always symmetric.

5 The Parity Fitness Landscape

Treat the landscape as a graph where the nodes are k length bit vectors. As in Section 4.2, if bit i is 1 this indicates the tree has an even number of D_i leafs. Nodes in the graph are directly connected (i.e. the corresponding trees are neighbours) if mutating one node gives the other in exactly one step. We deal with the 2^{k-1} nodes that are indirectly connected to the solution node.

Neighbours in the graph have exactly two bits different. Thus the graph consists of 2^{k-1} nodes each connected by $\frac{1}{2}k(k - 1)$ symmetric links and the probability of moving along each link is the same (i.e. $\frac{2}{k(k-1)}$). Only the single node with none of the k elements are even has fitness 1. All other nodes have zero fitness. Average fitness is $\bar{f} = 2^{1-k}$. The variance of fitness is $\frac{1}{2^{k-1}} \sum f_i^2 - \bar{f}^2 = \frac{1}{2^{k-1}} - \bar{f}^2 = 2^{1-k} - 2^{2-2k}$. So the standard deviation, $\sigma_f = \sqrt{2^{1-k} - 2^{2-2k}} \approx 2^{-(k-1)/2}$.

5.1 Fitness Distance Correlation

Jones and Forrest state that the fitness distance correlation based on random sampling of a needle in a haystack fitness function will be near zero [5, page 186]. We confirm this by giving values based on analysing the whole space and thus avoiding noise introduced by random sampling.

Jones and Forrest [5, page 185] define the fitness distance correlation as $r = \text{Cov}(f, d) / \sigma_f \sigma_d$. Where the covariance between fitness and distance (d) to the global optimum is: $\text{Cov}(f, d) = \frac{1}{2^{k-1}} \sum_{i=0}^{2^{k-1}-1} (f_i - \bar{f})(d_i - \bar{d})$ and σ_f is the standard deviation of f (calculated in the previous section) and similarly σ_d is the standard deviation of the number of two bit flips to the origin (i.e. distance to the global optimum).

$$\begin{aligned}
\text{Cov}(f, d) &= \frac{1}{2^{k-1}} \sum_{i=0}^{2^{k-1}-1} (f_i - \bar{f})(d_i - \bar{d}) \\
&= \frac{1}{2^{k-1}} \left((1 - \bar{f})(d_0 - \bar{d}) + \sum_{i=1}^{2^{k-1}-1} -\bar{f}(d_i - \bar{d}) \right) \\
&= \frac{1}{2^{k-1}} \left((1 - \bar{f})(d_0 - \bar{d}) - \bar{f} \sum_{i=1}^{2^{k-1}-1} d_i - \bar{d} \right) \\
&= \frac{1}{2^{k-1}} \left((1 - \bar{f})(-\bar{d}) + \bar{f}(d_0 - \bar{d}) - \bar{f} \sum_{i=0}^{2^{k-1}-1} d_i - \bar{d} \right) \\
&= \frac{1}{2^{k-1}} \left((1 - \bar{f})(-\bar{d}) - \bar{f} \bar{d} \right) \\
&= \frac{-\bar{d}}{2^{k-1}} \left((1 - \bar{f}) + \bar{f} \right) \\
&= \frac{-\bar{d}}{2^{k-1}}
\end{aligned}$$

The distance from the optimum is $bc/2$ where bc is the number of 1s (the bit count). If we round the division by two upwards we can just consider the bit count of the lower $k-1$ bits. (However for large k , $\lceil bc/2 \rceil$ can be approximated by $bc/2$.)

The number of points in the search space with identical bit count is C_{bc}^{k-1} . C_{bc}^{k-1} are coefficients of the binomial distribution. The mean of the binomial distribution is np and the variance is $np(1-p)$. So the mean divided by the standard deviation is $np / \sqrt{np(1-p)} = \sqrt{np/(1-p)}$. Here $p = 1/2$ and $n = k-1$. Assuming $\lceil bc/2 \rceil \approx bc/2$, the mean distance divided by the standard deviation of the distance $\bar{d} / \sigma_d \approx n^{1/2} = \sqrt{k-1}$.

$$\begin{aligned}
r &= \frac{\text{Cov}(f, d)}{\sigma_f \sigma_d} \approx \frac{-\bar{d}}{2^{k-1}} \frac{1}{2^{-(k-1)/2} \sigma_d} \\
&= \frac{-1}{2^{(k-1)/2}} \frac{\bar{d}}{\sigma_d} \\
r &\approx \frac{-\sqrt{k-1}}{2^{(k-1)/2}}
\end{aligned}$$

Figure 3 plots this approximation and the exact fitness distance correlation. It shows even for quite modest order, the actual correlation coefficient converges to this large order approximation.

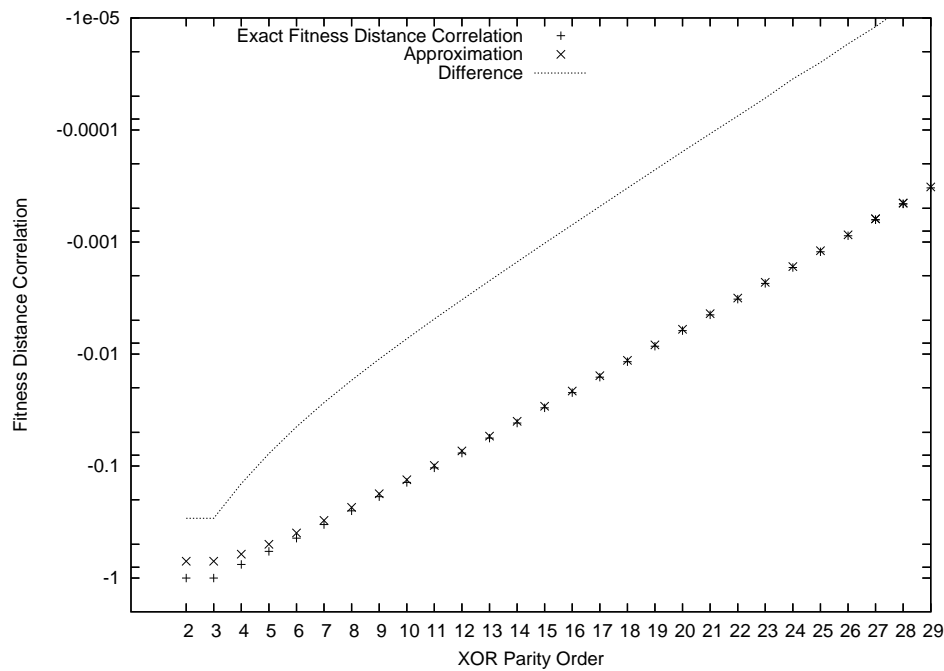


Figure 3: Fitness Distance Correlation for Needle in a Haystack with two bit flip mutation. Note log scale.

5.2 Walsh Analysis of Fitness Autocorrelation

Jones [5] defined fitness distance correlation in terms of distance from the global optimum. This has a number of problems. Usually one needs to know the location of the best solution in the search space and it is assumed that there is only one solution with the best fitness value. Fitness autocorrelation, for example along a random walk, addresses these problems and can be used as a measure of landscape smoothness and indicator of how easy a problem might be to solve. This section provides an informal argument that fitness correlation between neighbours falls rapidly with distance. It could be argued that in the case of a needle in a haystack landscape, such as parity, the chance of encountering the needle is sufficiently remote that fitness along a random walk will always be zero in practise. Nevertheless the following arguments can be applied to any landscape where the fitness landscapes created by the Walsh functions are elementary landscapes. This includes both 2-bit flip and 1-bit flip mutation.

We shall see in Section 5.10 that GP Parity can be represented as a linear sum of elementary landscapes. Dimova *et al.* [4] proved that fitness autocorrelation for a random walk on an elementary landscapes falls exponentially. Since the autocorrelation of all the components of parity fall monotonically, the correlation of parity itself must also fall monotonically with distance. (With increasing distance the actual value will be dominated by the slowest changing exponential term.)

Using Walsh analysis, any discrete fitness landscape can be represented as a sum of Walsh coefficients. If the Walsh basis functions are elementary landscapes (as is the case with parity) and fitness correlation falls rapidly in one step it will continue to fall towards zero for all larger distances.

For simplicity let the mean fitness be zero. Then fitness auto-correlation is essentially given by:

$$\sum_i \sum_{id} f_i f_{id}$$

where \sum_i is the sum over the whole search space and \sum_{id} is the sum over all neighbours of i which are d steps from i . For illustration assume $f = w + v$ where w and v are two Walsh basis elementary landscapes.

$$\begin{aligned} \sum_i \sum_{id} f_i f_{id} &= \sum_i \sum_{id} (w_i + v_i)(w_{id} + v_{id}) \\ &= \sum_i \sum_{id} w_i w_{id} + v_i w_{id} + w_i v_{id} + v_i v_{id} \end{aligned}$$

Now the first and last terms are simply the auto-correlation of w and v which we know fall exponentially.

$$\begin{aligned} \sum_i \sum_{id} v_i w_{id} &= \sum_i \sum_{id} v_i (w_i - (w_i - w_{id})) \\ &= \sum_i D v_i w_i - \sum_i \sum_{id} v_i (w_i - w_{id}) \end{aligned}$$

where D is the number of neighbours separated by d . Now the first term is zero, since w and v are orthogonal Walsh basis functions.

$$\sum_i v_i \sum_{id} (w_i - w_{id}) \leq v_{\max} \sum_i \sum_{id} (w_i - w_{id})$$

When the separation distance becomes large wrt w 's order $\sum_i \sum_{id} w_i - w_{id}$ will become small. Hence the autocorrelation between the sum of two Walsh functions falls rapidly with distance. This argument can be generalised to summing more than two elementary landscapes and holds for any fitness landscape where the Walsh basis functions are elementary landscapes.

As well as parity, there are interesting combinatorial problems where it is known the Walsh basis functions are elementary landscapes. (E.g. max-3-sat [20].) Fitness distance correlation will fall rapidly with distance for all of them.

5.3 Laplacian of the Parity Landscape Graph

Form the Laplacian Δ matrix as a real square matrix whose rows and columns correspond to the 2^{k-1} nodes in the parity landscape. Every element is zero except the off diagonal terms corresponding to an edge in the graph and the diagonal. Since the graph's edges are bidirectional and each have the same probability, the graph is symmetric and we can make all the non-zero off-diagonal terms be -1. The diagonal elements are the number of edges connected to the corresponding node in the graph. This is $\frac{1}{2}k(k-1)$ for all nodes. I.e. every node in the graph has the same degree. ($d = \frac{1}{2}k(k-1)$.) Thus every row (and every column) sums to zero.

$\Delta = \frac{1}{2}k(k-1)I - A$ where A is the graph's adjacency matrix.

$$\Delta_3 = \begin{array}{c} \begin{array}{cccc} & 000 & 101 & 110 & 011 \\ 000 & \left| \begin{array}{cccc} 3 & -1 & -1 & -1 \\ -1 & 3 & -1 & -1 \\ -1 & -1 & 3 & -1 \\ -1 & -1 & -1 & 3 \end{array} \right. \end{array} \end{array}$$

$$\Delta_4 = \begin{array}{c} \begin{array}{cccccccc} & 0000 & 1001 & 1010 & 0011 & 1100 & 0101 & 0110 & 1111 \\ 0000 & \left| \begin{array}{cccccccc} 6 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & 0 \\ -1 & 6 & -1 & -1 & -1 & -1 & -1 & 0 & -1 \\ -1 & -1 & 6 & -1 & -1 & 0 & -1 & -1 & -1 \\ -1 & -1 & -1 & 6 & 0 & -1 & -1 & -1 & -1 \\ -1 & -1 & -1 & 0 & 6 & -1 & -1 & -1 & -1 \\ -1 & -1 & 0 & -1 & -1 & 6 & -1 & -1 & -1 \\ -1 & 0 & -1 & -1 & -1 & -1 & 6 & -1 & -1 \\ 0 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & 6 \end{array} \right. \end{array} \end{array}$$

$$\Delta_5 = \begin{array}{c|cccccccccccccccc}
00000 & 10 & -1 & -1 & -1 & -1 & -1 & -1 & 0 & -1 & -1 & -1 & 0 & -1 & 0 & 0 & 0 \\
10001 & -1 & 10 & -1 & -1 & -1 & -1 & 0 & -1 & -1 & -1 & 0 & -1 & 0 & -1 & 0 & 0 \\
10010 & -1 & -1 & 10 & -1 & -1 & 0 & -1 & -1 & -1 & 0 & -1 & -1 & 0 & 0 & -1 & 0 \\
00011 & -1 & -1 & -1 & 10 & 0 & -1 & -1 & -1 & 0 & -1 & -1 & -1 & 0 & 0 & 0 & -1 \\
10100 & -1 & -1 & -1 & 0 & 10 & -1 & -1 & -1 & -1 & 0 & 0 & 0 & -1 & -1 & -1 & 0 \\
00101 & -1 & -1 & 0 & -1 & -1 & 10 & -1 & -1 & 0 & -1 & 0 & 0 & -1 & -1 & 0 & -1 \\
00110 & -1 & 0 & -1 & -1 & -1 & -1 & 10 & -1 & 0 & 0 & -1 & 0 & -1 & 0 & -1 & -1 \\
10111 & 0 & -1 & -1 & -1 & -1 & -1 & -1 & 10 & 0 & 0 & 0 & -1 & 0 & -1 & -1 & -1 \\
11000 & -1 & -1 & -1 & 0 & -1 & 0 & 0 & 0 & 10 & -1 & -1 & -1 & -1 & -1 & -1 & 0 \\
01001 & -1 & -1 & 0 & -1 & 0 & -1 & 0 & 0 & -1 & 10 & -1 & -1 & -1 & -1 & 0 & -1 \\
01010 & -1 & 0 & -1 & -1 & 0 & 0 & -1 & 0 & -1 & -1 & 10 & -1 & -1 & 0 & -1 & -1 \\
11011 & 0 & -1 & -1 & -1 & 0 & 0 & 0 & -1 & -1 & -1 & -1 & 10 & 0 & -1 & -1 & -1 \\
01100 & -1 & 0 & 0 & 0 & -1 & -1 & -1 & 0 & -1 & -1 & -1 & 0 & 10 & -1 & -1 & -1 \\
11101 & 0 & -1 & 0 & 0 & -1 & -1 & 0 & -1 & -1 & -1 & 0 & -1 & -1 & 10 & -1 & -1 \\
11110 & 0 & 0 & -1 & 0 & -1 & 0 & -1 & -1 & -1 & 0 & -1 & -1 & -1 & -1 & 10 & -1 \\
01111 & 0 & 0 & 0 & -1 & 0 & -1 & -1 & -1 & 0 & -1 & -1 & -1 & -1 & -1 & -1 & 10
\end{array}$$

For a landscape to be elementary its fitness function must have a special relationship with its move operator; it must obey the “wave equation” $\Delta f = \lambda(f - \bar{f})$ (see Section 2.1). I.e. treat the fitness function as a vector whose elements are the fitnesses of the corresponding nodes in landscape graph. Call this f . If the landscape is to be elementary f (up to an additive constant) must be an eigenvector of Δ with corresponding eigenvalue λ . For parity, f is a vector with 2^{k-1} elements all of which are zero except the first, which has the value 1. Thus Δf is simply the first column of Δ . The first element of the first column of Δ is $\frac{1}{2}k(k-1)$ and there are $\frac{1}{2}k(k-1)$ other elements whose values are -1. Thus the first column of Δ is not a simple scalar multiple of the fitness vector f and so f is not an eigenvector of Δ . Therefore the parity landscape is not elementary.

Section 5.5 and those following it will show not only is parity’s landscape not elementary but also it cannot be decomposed into a reasonably small number of elementary landscapes.

5.4 Elementary Needle in Haystack Landscape

We use the analysis from the previous section to construct another needle in a haystack (NIH) problem which is elementary. Without loss of generality we can keep the solution at all zeros and give it fitness one (all other points have zero fitness). Thus the NIH has the same fitness function as parity. There are 2^l points in the search space, so $\bar{f} = 2^{-l}$. For an NIH landscape to be elementary it must still obey the “wave equation” $\Delta f = \lambda(f - \bar{f})$ and Δf is again the first column of Δ . An elementary NIH Laplacian is:

$$\Delta_{\text{NIH}} = \begin{array}{cccccc}
& 2^l-1 & -1 & -1 & \dots & -1 \\
& -1 & 2^l-1 & -1 & \dots & -1 \\
\Delta_{\text{NIH}} = & -1 & -1 & 2^l-1 & \dots & -1 \\
& \vdots & \vdots & \vdots & \ddots & \vdots \\
& -1 & -1 & -1 & \dots & 2^l-1
\end{array}$$

which gives $\lambda = 2^l$. That is we can construct an elementary needle in a haystack landscape with a mutation operator which can move to any point in the search landscape in one step and all such points are equally likely.

5.5 Recursive Construction of Parity's Landscape

This section will show that the Laplacian of the connectivity graph for k order parity (Δ_k) can be recursively created from two Laplacian for $(k-1)$ parity and two Laplacians for the $k-2$ Hamming graph (H_{k-2}). I.e:

$$\Delta_k = \left(\begin{array}{c|c} (k-1)I + \Delta_{k-1} & H_{k-2} - (k-1)I \\ \hline H_{k-2} - (k-1)I & (k-1)I + \Delta_{k-1} \end{array} \right) \quad (2)$$

Where I is the identity matrix (of the appropriate dimensions). $(k-1)I$ is included in the two on-diagonal sub-matrices so that all the diagonal elements are $(k-1) + \frac{1}{2}(k-1)(k-2) = \frac{1}{2}k(k-1)$ as required. Since H (like Δ) is a graph Laplacian, its rows sum to zero. By subtracting $(k-1)I$ from the off diagonal matrices we ensure all the rows in $(k-1)I + \Delta_{k-1} | H_{k-2} - (k-1)I$ also sum to zero. (When we look in detail at H_{k-2} we will see that $-(k-1)I$ is exactly the adjustment needed for the Hamming cube adjacency matrix.)

Label the rows and columns of Δ_k with 2^{k-1} integers. These are given by the first 2^{k-1} integers starting from zero plus a k^{th} bit. This most significant bit is set (or cleared) to ensure the number of bits set in the row label is even. Order the rows/columns by the lower $k-1$ bits, ignoring the top (parity) bit. Elements of Δ are -1 if there are exactly two bits different in the row and column labels. Diagonal elements are $\frac{1}{2}k(k-1)$ and all other elements are zero.

Divide Δ into four equal sub-matrices. Half of Δ is made of the rows whose $k-1^{th}$ bit is zero. The other half is made of the rows whose $k-1^{th}$ bit is one. (Similarly for the columns). In the two on-diagonal sub-matrices the $k-1^{th}$ bit in both the rows and the columns is the same.

If an element of Δ in the on-diagonal sub-matrices is -1, then this means that there are exactly two bits that differ in its row and column labels but the differing bits cannot include the $k-1^{th}$ bit. If we remove the $k-1^{th}$ bit, this is the condition for Δ_{k-1} to be -1. (We have to tidy up the labels by not just removing the $k-1^{th}$ bit but also the k^{th} bits and recalculating the top, parity, bit for the $k-2$ bits.) As mentioned above, Δ_{k-1} has diagonal elements of $\frac{1}{2}(k-1)(k-2)$ so $k-1$ has to be added to them to convert the diagonal elements of Δ_{k-1} to those for Δ_k .

In the off diagonal sub-matrices of Δ , either the $k-1^{th}$ bit of the row's label is zero and the $k-1^{th}$ bit of the column's label is one or vice versa. I.e. the row and column's label already differs in one bit. If an element of Δ in the off-diagonal sub-matrices is -1, then, excluding the $k-1^{th}$ and k^{th} bits, its row and column labels must differ by exactly one bit. (Alternatively the $k-1^{th}$ and k^{th} bits both differ and the other $k-2$ bits are the same.) Ignoring the top two bits for a moment, we see that Δ being -1 in the off diagonals is exactly the same as the Hamming distance between two $k-2$ bit strings being one. That is, if H_{k-2} is -1 so too are the corresponding off diagonal elements in Δ . The only other non-zero elements of H are on the diagonal.

Since a k -bit string has k neighbours which differ by exactly one bit, the diagonal elements of H are k . I.e., the diagonal elements of H_{k-2} are $k-2$. These elements correspond to the lower $k-2$ bits of the row and column labels of Δ being the same. In Δ these elements are -1 (rather than $k-2$) since the two top bits can be simultaneously changed without changing the lower $k-2$ bits. Subtracting $k-1$ from the diagonal elements of H_{k-2} (i.e. from $k-2$) gives -1. Which is the value of the corresponding element in Δ_k . Hence we have proved Equation 2.

Δ_{k-1} can be defined in terms of Δ_{k-2} and H_{k-3} and so on. The base cases are: Δ_2 and H_1 which are both $\begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix}$.

5.6 Eigen Analysis of Parity's Landscape

Using the recursive decomposition of Δ given in the previous section, we shall show if e_{k-1} is an eigenvector of Δ_{k-1} then $e_{k-1}|e_{k-1}$ and $e_{k-1}| - e_{k-1}$ are eigenvectors of Δ_k . It turns out the Walsh functions [19] are eigenvectors of both the Laplacian of the Hamming neighbourhood H and of the that of the Parity neighbourhood Δ . The eigenvalues of Δ will be given at the end of this section.

$$e_k \Delta_k = (e_{k-1}| \pm e_{k-1}) \left(\begin{array}{c|c} (k-1)I + \Delta_{k-1} & H_{k-2} - (k-1)I \\ \hline H_{k-2} - (k-1)I & (k-1)I + \Delta_{k-1} \end{array} \right) = \left(\begin{array}{cc} (k-1)e_{k-1} + \lambda e_{k-1} & \pm e_{k-1} H_{k-2} \mp (k-1)e_{k-1} \\ e_{k-1} H_{k-2} - (k-1)e_{k-1} \pm & (k-1)e_{k-1} \pm \lambda e_{k-1} \end{array} \right) \quad (3)$$

Start with the base case ($k=2$): $\begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix}$ has eigenvectors $e_{2+}=(1,1)$ (eigenvalue 0) and $e_{2-}=(1,-1)$ (eigenvalue 2). Notice that these are eigenvectors of both Δ_2 and H_1 . So Equation 3 for $k=3$ becomes

$$(e_2| \pm e_2) \Delta_3 = \left(\begin{array}{cc} (k-1)e_{k-1} + \lambda e_{k-1} & \pm \lambda e_{k-1} \mp (k-1)e_{k-1} \\ \lambda e_{k-1} - (k-1)e_{k-1} & \pm (k-1)e_{k-1} \pm \lambda e_{k-1} \end{array} \right) = \left(\begin{array}{cc} 2e_2 + \lambda e_2 & \pm \lambda e_2 \mp 2e_2 \\ \lambda e_2 - 2e_2 & \pm 2e_2 \pm \lambda e_2 \end{array} \right) = \left(\begin{array}{c} 2\lambda e_2 \\ 2\lambda e_2 \end{array} \right) \text{ and } \left(\begin{array}{c} 4e_2 \\ -4e_2 \end{array} \right)$$

That is $(e_2|e_2)$ are eigenvectors of Δ_3 (with eigenvalue 2λ) and so too are $(e_2| - e_2)$ (with eigenvalue 4). So Δ_3 has four eigenvectors:

$$\begin{aligned} e_{3++} &= (1, 1, 1, 1) \\ e_{3+-} &= (1, 1, -1, -1) \\ e_{3-+} &= (1, -1, 1, -1) \\ e_{3--} &= (1, -1, -1, 1) \end{aligned}$$

with corresponding eigenvalues: $\lambda_{3++} = 2\lambda_{2+} = 0$, $\lambda_{3+-} = 2\lambda_{2-} = 4$, $\lambda_{3-+} = 4$ and $\lambda_{3--} = 4$. We can see that the four eigenvectors form an orthonormal set and so this is a complete eigen description of Δ_3 . The eigenvectors of Δ_3 , e_3 are the Walsh basis (on 2 bits).

It can be shown that the Walsh basis (on k bits) are eigenvectors of the Hamming cube H_k , with eigenvalues $2i$ with multiplicity C_i^k for $0 \leq i \leq k$.¹ Therefore the eigenvectors of Δ_3 are also eigenvectors of H_2 (albeit with different eigenvalues). We have e_3 are the Walsh basis therefore $e_3|e_3$ and $e_3| - e_3$ are the Walsh basis (3 bits). Let this hold for any higher order. I.e. $e_{k-1}|e_{k-1}$ and $e_{k-1}| - e_{k-1}$ are the Walsh basis on $k-1$ bits. Then Equation 3 becomes:

$$\left(\begin{array}{c|c} (k-1)e_{k-1} + \lambda_{k-1}e_{k-1} \pm 2ie_{k-1} \mp (k-1)e_{k-1} & \\ \hline 2ie_{k-1} - (k-1)e_{k-1} \pm (k-1)e_{k-1} \pm \lambda_{k-1}e_{k-1} & \end{array} \right) = \left(\begin{array}{c} (2i + \lambda_{k-1})e_{k-1} \\ (2i + \lambda_{k-1})e_{k-1} \end{array} \right) \text{ and } \left(\begin{array}{c} (2(k-1) - 2i + \lambda_{k-1})e_{k-1} \\ -(2(k-1) - 2i + \lambda_{k-1})e_{k-1} \end{array} \right)$$

Where $2i$ is an eigenvalue of H_{k-2} . So $0 \leq i \leq k-2$. $2i$ and λ_{k-1} are related since they are the eigenvalues of H_{k-2} and Δ_{k-1} for the same eigenvector e_{k-1} .

¹See, for example, Dr. Daniel Spielman's lecture notes (Lecture five, 16 September 2009 <http://www.cs.yale.edu/homes/spielman/561/lect05-09.pdf>) or [2].

That is $e_{k-1}|e_{k-1}$ are indeed eigenvectors of Δ_k (with eigenvalues $2i + \lambda_{k-1}$) and so too are $e_{k-1}| - e_{k-1}$ (with eigenvalues $2(k-1) - 2i + \lambda_{k-1}$). Notice since $e_{k-1}|\pm e_{k-1}$ are the Walsh basis, they form a complete orthogonal set of eigenvectors for Δ_k . The eigenvalues of Δ are not as elegant as those of the Hamming cube H but can be rapidly calculated, $O(k^2)$. The first values are given in Table 1. Notice there are rather fewer distinct values than for the Hamming cube. In Section 5.9 we will prove that Δ_k has $k/2+1$ distinct eigenvalues.

The multiplicities shown as subscripts in Table 1 are similar to Pascal’s triangle. (I.e. the eigenvalue multiplicities of the Hamming cube’s Laplacian.) Firstly they sum to 2^{k-1} in each row. Also, except for the largest eigenvalue, each multiplicity is the sum of the multiplicities immediately above and to the left in the previous row. The multiplicity of the largest eigenvalue depends upon whether k is odd or even. If k is even, the multiplicity of the largest eigenvalue is the same as that in the previous row. If k odd, it is the sum of the multiplicity in the previous row to the left plus *twice* the multiplicity of the largest eigenvalue for $k-1$.

5.7 Largest Eigenvalue of 2-bit Flip Graph Laplacian

Since Δ is a Laplacian and hence its rows always sum to zero, its smallest eigenvalue is always zero (with multiplicity one). The smallest non-zero eigenvalue corresponds to the lowest Walsh function and has the value $2k-2$. Saying in closed form which Walsh function has the largest eigenvalue is not straight forward. However we can 1) establish an upper bound on the largest eigenvalue and 2) give a stochastic estimate for large k .

5.7.1 Upper Bound on Largest Eigenvalue of 2-bit Flip Graph Laplacian

Depending upon the Walsh function chosen, an eigenvalue of Δ_k is either $\lambda_{H_{k-2}} + \lambda_{k-1}$ or $2(k-1) - \lambda_{H_{k-2}} + \lambda_{k-1}$. We know $\lambda_{H_{k-2}}$ cannot exceed $2k-4$ or be smaller than zero. Therefore the largest eigenvalue cannot exceed $2k-2 + \max(\lambda_{k-1})$. We know $\max(\lambda_2) = 2$. So (for $k > 2$) $\lambda_k \leq \sum_{i=3}^{i=k} 2i - 2 + 2$. Rearranging gives $\lambda_k \leq k(k+1) - 6$. However it appears the actual value is $\lceil (k-1)(k+1)/2 \rceil$.

Table 1: Eigenvalues of Laplacian of Parity’s landscape graph. The subscripts are the multiplicity of the Δ eigenvalues. Last column is the number of distinct eigenvalues. The total number of eigenvalues and eigenvectors is 2^{k-1} .

$k = 2$	0_1	2_1					2			
$k = 3$	0_1	4_3					2			
$k = 4$	0_1	6_4	8_3				3			
$k = 5$	0_1	8_5	12_{10}				3			
$k = 6$	0_1	10_6	16_{15}	18_{10}			4			
$k = 7$	0_1	12_7	20_{21}	24_{35}			4			
$k = 8$	0_1	14_8	24_{28}	30_{56}	32_{35}			5		
$k = 9$	0_1	16_9	28_{36}	36_{84}	40_{126}			5		
$k = 10$	0_1	18_{10}	32_{45}	42_{120}	48_{210}	50_{126}			6	
$k = 11$	0_1	20_{11}	36_{55}	48_{165}	56_{330}	60_{462}			6	
\vdots										
$k = 17$	0_1	32_{17}	60_{136}	84_{680}	104_{2380}	\cdots	144_{24310}			9

5.7.2 Long Bit String Estimate of the Largest Eigenvalue of the 2-bit Flip Graph Laplacian

For any unit vector u the length of Δu will be no bigger than the largest eigenvalue. Choose a random direction. I.e. let v be a vector of 2^{k-1} components each of which is either +1 or -1, chosen uniformly at random. $|v|^2 = 2^{k-1}$ hence $|v| = 2^{(k-1)/2}$ (Eventually we will normalise v to be of unit length by dividing by $|v|$.) The first element of Δv is typical of them all.

$$\Delta v(1) = \pm \frac{1}{2}k(k-1) \underbrace{\pm 1 \pm \dots \pm 1}_{\frac{1}{2}k(k-1) \text{ terms with random signs}}$$

The square of the length is given by summing the squares of each component in the usual way. Since the elements of v were randomly chosen, each of the elements of Δv are independent and identically distributed and therefore $|\Delta v(i)|^2$ are also i.i.d. Thus the expected value of $|\Delta v|^2$ is $2^{k-1} \times$ the expected value of any of them, e.g. the first $|\Delta v(1)|^2$. The expected length of $\Delta v(1)$ is approximately $\frac{1}{2}k(k-1)$. (The random signs ensure on average the following terms come to near zero and for large k the sum is dominated by the first (largest) term.) Thus the expected value of $|\Delta v|^2$ is $|\frac{1}{2}k(k-1)|^2 2^{k-1}$ and that of $|\Delta v| = \frac{1}{2}k(k-1)2^{(k-1)/2}$. Taking the ratio $\frac{|\Delta v|}{|v|}$ gives $\frac{1}{2}k(k-1)$ as an upper bound on all the eigenvalues. For large k this will become a tight bound on the largest eigenvalue.

Note the exact bound appears to be within a factor of two of the apparent value, whereas the stochastic upper bound appears to be increasingly tight as k increases. Finally note the eigenvalues of parity's Δ are a factor of k bigger than those of the Hamming cube.

5.8 Elementary Landscape Roughness and the Eigenvalues of the Graph Laplacian

Recall $d = \frac{1}{2}k(k-1)$ so the stochastic bound is needed to ensure $\frac{\lambda}{d} \leq 1$ for all cases (cf. Equation 1). Higher order Walsh functions are considered to be more rugged, since their sign changes more often than lower order Walsh functions. (Rothlauf points out [14, p27] that Walsh order is not a universal indicator of problem difficulty.) Elementary landscapes generated by higher Walsh basis functions have higher eigenvalues than those corresponding to lower order Walsh functions. Using Equation 1, we can see the higher an elementary landscape's eigenvalue the further each point is from the average of its neighbours. Thus we can view $\frac{\lambda}{d}$ as another measure of landscape ruggedness. The larger it is, the less each point tells us about the (average) fitness of its neighbours.

In other regular elementary landscapes eigenvalues can exceed the number of neighbours. I.e. $\frac{\lambda}{d}$ can exceed 1. For example in the Hamming cube the largest eigenvalue of H_k is $2k$ and each node has k neighbours. (So $0 \leq \lambda_H/d_H \leq 2$.) Whitley *et al.* [21, p589] equates $\frac{\lambda}{d} > 1$ with rugged elementary landscapes. Whereas they suggests if $0 \leq \frac{\lambda}{d} \leq 1$ the elementary landscape is smooth. ([21] uses a constant rather than referring to λ as an eigenvalue.)

We get the same conclusion if we use Whitley's [20] component based model of elementary landscapes. This treats their fitness as being composed of components which are added to and removed from the current trial solution as the search process moves from a point to one of its neighbours. [20, p383] gives $\frac{\lambda}{d} = p_1 + p_2$. Where p_1 is the proportion of components of $f(x)$ that change when we move away from x . p_2 is the proportion of components not included in $f(x)$ that change when we move away from x . Thus we should expect a small value of $p_1 + p_2$ to give a smooth landscape. When $p_1 + p_2$ approaches one, most of the components are being changed at each move, so we expect a more rugged landscape, in keeping with the previous paragraph.

5.9 Number of Distinct Eigenvalues and Parity's Graph Diameter

A graph's diameter is the maximum distance between any two nodes in the graph. (Where the distance is the smallest number of edges that have to be traversed to go between the nodes.) Parity is symmetric so the longest distance between any pair of nodes, is the same as the longest distance between the origin and any node. This is the minimum number of pairs of bit flips between $k-1$ zeros and the binary string of the target node. This is simply $bc(\text{target})/2$, where bc is the number of 1s (the bit count). Obviously the worst case is when all bits are one. Hence Δ 's graph diameter is $k/2$. Whereas the diameter of the Hamming cube is $k-1$.

Reeves [12, page 598] says the number of distinct eigenvalues is the graph diameter plus one. We can see, from the last column of Table 1, that the number of distinct eigenvalues is indeed $k/2+1$.

5.10 Number of Distinct Elementary Landscapes and Walsh Analysis

Since the eigenvectors form an orthogonal set, any vector, including the fitness f vector, can be represented by its components projected onto the eigenvectors. Remembering Section 5.3, each eigenvector of Δ represents an elementary landscape, so fitness can be represented as a sum of elementary landscapes. One for each eigenvector where it has a non-zero projection. As the Walsh basis functions are eigenvectors of Δ projecting the fitness f vector onto these eigenvectors is equivalent to the Walsh analysis of f . Since f is 1 followed by $2^{k-1} - 1$ zeros it has 2^{k-1} non-zero Walsh coefficients. That is all its Walsh coefficients are non-zero (actually they are all equal to $2^{-(k-1)}$). However eigenvectors with the same eigenvalue form sub-spaces where any linear combination of these eigenvectors is also an eigenvector. Hence any vector can be represented as a sum of eigenvectors, one for each subspace where it has non-zero projection. We know that Parity's fitness vector has non-zero projection into each subspace so all $k/2 + 1$ subspaces must be used. That is, Parity's fitness function can be expressed as $k/2 + 1$ elementary fitness landscapes. Indeed $k/2 + 1$ is an upper limit on the number of elementary landscapes needed to represent any fitness function (when we use only a 2-bit flip mutation operator).

6 Comparison with Real GP

6.1 NIH Non-Elementary Landscape

To verify GP does behave similarly to our model, we first ran GP to show it treats EQ-parity as a needle in a haystack problem and to investigate the distribution of jump sizes in a real GP. In the first group of genetic programming runs, GP was run with Koza's 16-even parity fitness function and only EQ in the function set (details given in Table 2).

TinyGP uses the "grow" method [11] to create the initial random programs. Since there are four times as many terminals as functions, despite a large depth limit (8), the grow method produces populations that

Table 2: TinyGP Parameters for 16 even parity

Function:	EQ
Terminals:	D_0, \dots, D_{15}
Fitness:	Number of correct answers on all 2^{16} test cases. However (see Section 3) only fitness 32768 and 65536 are possible.
Selection:	Steady state. 2 members tournaments.
Population:	65 536
Initial pop:	grow (max depth 8),
Parameters:	80% subtree crossover, 20% point mutation (p_m 0.05). Crossover and mutation points are chosen uniformly (i.e. without a function bias [6]) No size limit.
Termination:	100 generations

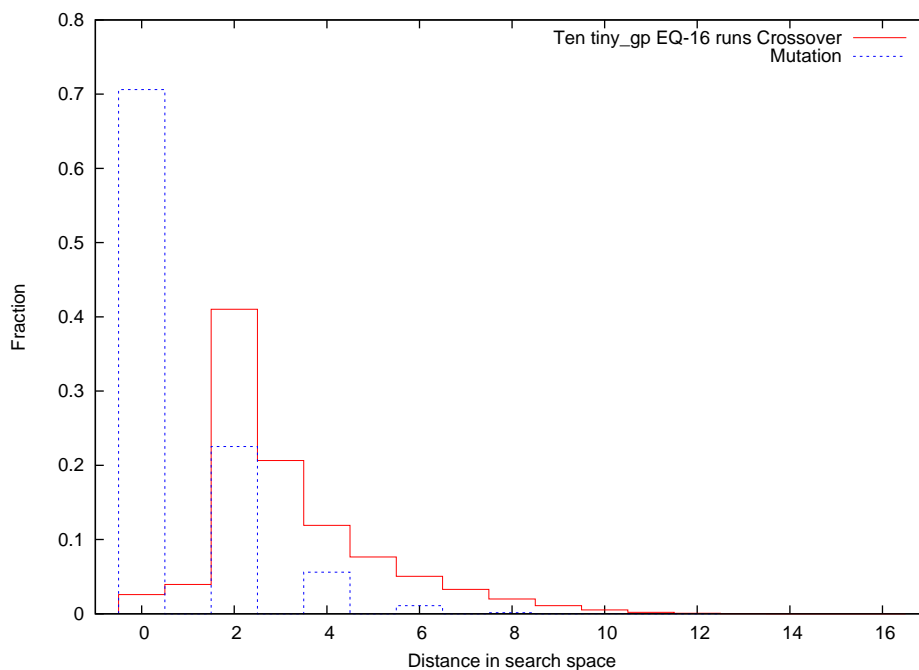


Figure 4: Distribution of jumps in the search space sizes caused by crossover and mutation in ten TinyGP 16-EQ parity runs.

consist mostly of programs that are too small (mean 17) to solve the 16-EQ parity problem (minimum solution size 31). Note since an NIH landscape does not provide fitness guidance, the population evolves as though there was no fitness, hence there is no bloat and, excluding drift, on average programs do not change size [7]. Therefore there remains a substantial part of the population which is simply too small to solve the problem. This increases the search time. Also TinyGP does not ensure children are different from their parents. This also increases the search time. And so a large number of programs need to be run before finding a solution. The first hitting time is 2 400 000 (mean of 10 runs, standard deviation 2 200 000). Also the distribution shows signs of being geometric as expected of an NIH problem.

Although the programs are rather smaller than assumed in Section 4.3, if we exclude mutations which make no difference (distance 0) the distribution of jump sizes for TinyGP mutation (Figure 4 dashed line) is somewhat similar to that predicted in Section 4.3. That is, most mutations *moves* cause the child to be exactly two bits different in our search space. The number of 4, 6, 8 etc. moves falls rapidly. (As an anti-bloat mechanism TinyGP uses point mutation with a fixed probability of mutation per program element [11]. Thus not only can zero elements be changed but also a mutant child may differ in multiple places from its mother.) Notice, as expected, mutation only causes jumps by multiples of 2.

Figure 4 also shows the distribution of jump sizes caused by crossover. Note since subtree crossover can change program sizes, it can make arbitrary jump sizes (including odd sized jumps). However the most popular (mode) jump size is two, as assumed for the mutation only model presented in Section 4.3,

Most crossovers take place between trees which were themselves created by crossovers. (Point mutation does not change tree size or shape and so need not be considered.) It may be this that gives the distribution (excluding 0 and 1) its pleasing Zipf like tail. (Falling by about 60% per unit increase in step size.) Note, although the distribution of tree sizes would be expected to rapidly converge to a Lagrange distribution [3], the distribution in Figure 4 refers to jumps in our bit orientated semantic search space. The mapping between it and that of the GP binary trees is not straightforward and we have not attempted to prove a mathematical relationship between random crossover of Lagrange distributed trees and the jumps in our space to accompany the experimental data given in Figure 4. Figure 4 shows there is some truth in our simplification that all moves are of size two. However it shows the full situation is more complex.

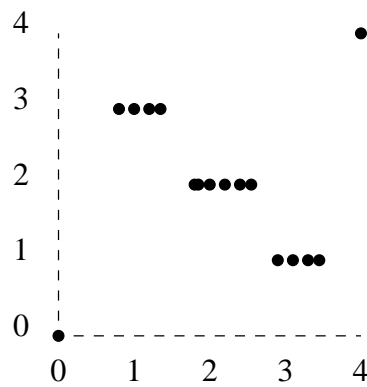


Figure 5: Schematic of the integer valued 4 bit unitation fitness function formed by adding four third order 4-bit Walsh functions. Fitness plotted vertically. Unitation horizontally. Spots indicate the fitness of 16 possible bit values. The full fitness function is formed by concatenating four of these. It has a single global optimum with fitness 16.

6.2 1 Max Elementary GP Landscape

A further ten runs were made where the fitness function is given by the Hamming distance from the optimum in our bit-orientated semantic landscape. (Note in GP terms this is cheating, since we allow the fitness function to look inside the program, rather than simply being based on what the program does when it is run.) The problem is now like ones max and so considerably easier than the parity problem used in the previous section. n^{th} order one-max can be shown to be an elementary landscape by either noting 1) its fitness is composed of n components in Whitley's sense [18] or 2) by noting it can be decomposed into the n first order Walsh coefficients all of which are eigenvectors of Δ with the same eigenvalue and hence any linear combination of them (such as onemax) is also an eigenvector of Δ and therefore onemax is elementary. (The non-zero zeroth Walsh coefficient essentially gives \bar{f} .)

TinyGP was run with a much reduced population size (128). Except for this and the fitness function, the same parameters (cf. Table 2) were used. 85% of runs succeeded in solving EQ-16 on this more friendly landscape. The mean first hitting time was 1020 (standard deviation 430). I.e. the second fitness landscape is about 2000 times easier for GP even though their search spaces are identical.

6.3 4 Trap-like Third Order GP Landscape

A further ten runs were made where the fitness function is given by summing four trap-like functions. Each trap-like function is defined on a group of four non-overlapping bits from the total of 16 in our bit-orientated semantic landscape. Where each four bit function is created by adding together all the third order Walsh coefficients. (See Figure 5. Again in GP terms this is cheating.)

If we retain Section 4's assumption that search proceeds by flipping exactly two bits, then the new fitness landscape is elementary. To show this we start with the four bit (i.e. 16×16) Laplacian Δ_5 and note the fitness function is composed of third order Walsh functions and repeat the argument given on page 11 but for this special case. We have already shown the third order Walsh functions are eigenvectors of the 16×16 Laplacian and have the same eigenvalue (12). The 32×32 Laplacian is formed of four 16×16 matrices as described in Equation 2. Concatenating a 16-bit Walsh function with itself and multiplying by the 32×32 Laplacian gives a vector of 32 elements, whose two 16 elements halves are both equal to a multiple of the original 16-bit Walsh function. I.e. the 32 element vector is an eigenvector of the 32×32 Laplacian. The eigenvalue is λ_{H_3} plus that of the original 16-bit Walsh function. Where λ_{H_3} is the eigenvalue of the third order Walsh function applied to the Laplacian of the Hamming cube (rather than two bit flip parity). $\lambda_{H_3} = 2 \times 3 = 6$. I.e. the concatenated 32 bit third order Walsh functions have eigenvalues $6 + 12=18$. Note this is true of all four original third order Walsh functions and therefore they all have the

same eigenvalue. Therefore the 32 element vector formed by adding them together is also an eigenvector with eigenvalue 18.

We keep doubling using this procedure until we form a vector with 2^{16} elements. It will be an eigenvector of the $2^{16} \times 2^{16}$ Laplacian for the 16-EQ search space; with eigenvalue $6 \times (16 - 4) + 12 = 84$. Notice this vector repeats the 16 values of the third order Walsh function 4096 times, corresponding exactly to the repeat of the lower 4 bits of our trap-like fitness function. Since our landscape is symmetric we can rotate the labels by 4, 8 and 12 bits to show that the three other components of the fitness function are also eigenvectors with the same eigenvalue. Since the combined fitness function is the sum of four eigenvectors with the same eigenvalue, it too is an eigenvector of the $2^{16} \times 2^{16}$ Laplacian (with eigenvalue 84.) Therefore our trap-like fitness landscape is elementary.

The problem is now intermediate between parity and onemax. Therefore TinyGP was run with an intermediate population size (1024). Except for this and the fitness function, the same parameters (cf. Table 2) were used. Nine out of ten runs succeeded. The mean first hitting time was 18 000 (standard deviation 9 000). I.e. this third fitness landscape is about 130 times easier for GP than the first (and about 17 times harder than the second) even though all three search spaces are identical.

7 Conclusions

We have analysed our [8] genetic programming parity fitness landscape. As well as proving it is not elementary we have calculated its mean fitness and fitness variance and its fitness distance correlation. Also we have argued that existing results on fitness autocorrelation along random walks in elementary landscapes can be greatly extended.

We have given an elementary needle in a haystack fitness landscape.

We have given a complete eigen analysis of the search space formed by using two bit flip mutation. Showing the eigenvectors of its graph Laplacian are the same as those of the single bit flip Hamming cube (i.e. the Walsh basis functions). The eigenvalues can be rapidly calculated. They are somewhat similar but a factor of about k larger than those of the Hamming cube. Since the graph is connected the smallest eigenvalue is of course zero. The next is $2(k-1)$, the next $4(k-2)$, then $6(k-3)$ and so on. Table 1 gives the eigenvalues and their multiplicities for initial values of k . We provide two proofs (one a bound and the other a tight asymptotic limit) for the largest eigenvalue. The actual values appear to be $\lceil (k-1)(k+1)/2 \rceil$ but we have not proved this. The number of distinct eigenvalues is $k/2+1$ and so the separation between eigenvalues is about $2k$. For large k , multiplying the Laplacian by a unit vector pointed in almost all directions will increase its length by about the number of neighbours of each node $\frac{1}{2}k(k-1)$. It may be possible to generalise this to landscapes where neighbours differ by n bits, with $n > 2$ [17].

The number of distinct eigenvalues is $k/2 + 1$ and total number of eigenvectors is 2^{k-1} . Therefore there is at least one subspace with at least $2^k/(k+2)$ eigenvectors. Any linear combination of these is also an eigenvector and each of these corresponds to an elementary landscape. Restricting ourselves to coefficients 0 and 1 means this subspace alone contains at least $2^{2^k/(k+2)}$ 2-bit flip elementary landscapes. (This is a lower bound, the total number of 2-bit flip elementary landscapes is much bigger.)

We compared our simplified [8] genetic programming parity with experiment and showed it does indeed behave as a needle in a haystack. We have run GP on two elementary fitness landscapes of the same size but different fitness functions, and shown, as expected, they behave differently. I.e. elementary landscapes, with identical sizes and connectivity, can represent problems of very different difficulty. Results so far have been in keeping with the ruggedness measure (Section 5.8). However given non-universal results on other proposed indicators of problem hardness (e.g. order of non-zero Walsh coefficients) we cannot be confident of its utility.

Acknowledgments

I would like to thank Daniel Spielman, Gylson Thomas for fhtseq.m, Riccardo Poli for TinyGP and Gabriel Peyr for fwt.m. Also many participants at Dagstuhl Seminar 10361 for helpful discussions, including Jon Rowe, Francisco Chicano, Andrew Sutton and Darrell Whitley, and FOGA 2011 anonymous reviewers.

References

- [1] T. Bäck. *Evolutionary Algorithms in Theory and Practice: Evolution Strategies, Evolutionary Programming, Genetic Algorithms*. Oxford University Press, New York, 1996.
- [2] T. Biyikoglu, J. Leydold, and P. F. Stadler. *Laplacian Eigenvectors of Graphs: Perron-Frobenius and Faber-Krahn Type Theorems*, volume 1915 of *Lecture Notes in Mathematics*. Springer, 2007. On line.
- [3] S. Dignum and R. Poli. Generalisation of the limiting distribution of program sizes in tree-based genetic programming and analysis of its effects on bloat. In D. Thierens, H.-G. Beyer, J. Bongard, J. Branke, J. A. Clark, D. Cliff, C. B. Congdon, K. Deb, B. Doerr, T. Kovacs, S. Kumar, J. F. Miller, J. Moore, F. Neumann, M. Pelikan, R. Poli, K. Sastry, K. O. Stanley, T. Stutzle, R. A. Watson, and I. Wegener, editors, *GECCO '07: Proceedings of the 9th annual conference on Genetic and evolutionary computation*, volume 2, pages 1588–1595, London, 7-11 July 2007. ACM Press.
- [4] B. Dimova, J. W. Barnes, and E. Popova. Arbitrary elementary landscapes & AR(1) processes. *Applied Mathematics Letters*, 18(3):287–292, 2005.
- [5] T. Jones and S. Forrest. Fitness distance correlation as a measure of problem difficulty for genetic algorithms. In *Proceedings of the 6th International Conference on Genetic Algorithms, ICGA 1995*, pages 184–192. Morgan Kaufmann, 1995.
- [6] J. R. Koza. *Genetic Programming: On the Programming of Computers by Natural Selection*. MIT press, 1992.
- [7] W. B. Langdon. The evolution of size in variable length representations. In *1998 IEEE International Conference on Evolutionary Computation*, pages 633–638, Anchorage, Alaska, USA, 5-9 May 1998. IEEE Press.
- [8] W. B. Langdon. Scaling of program tree fitness spaces. *Evolutionary Computation*, 7(4):399–428, Winter 1999.
- [9] W. B. Langdon and R. Poli. *Foundations of Genetic Programming*. Springer-Verlag, 2002.
- [10] R. Poli, W. B. Langdon, and O. Holland. Extending particle swarm optimisation via genetic programming. In M. Keijzer, A. Tettamanzi, P. Collet, J. I. van Hemert, and M. Tomassini, editors, *Proceedings of the 8th European Conference on Genetic Programming*, volume 3447 of *Lecture Notes in Computer Science*, pages 291–300, Lausanne, Switzerland, 30 Mar. - 1 Apr. 2005. Springer.
- [11] R. Poli, W. B. Langdon, and N. F. McPhee. *A field guide to genetic programming*. Published via <http://lulu.com> and freely available at <http://www.gp-field-guide.org.uk>, 2008. (With contributions by J. R. Koza).
- [12] C. R. Reeves. Fitness landscapes. In E. K. Burke and G. Kendall, editors, *Search Methodologies: Introductory Tutorials in Optimization and Decision Support Techniques*, chapter 19, pages 587–610. Springer, 2005.
- [13] C. R. Reeves and J. E. Rowe. *Genetic Algorithms—Principles and Perspectives: A Guide to GA Theory*. Kluwer Academic Publishers, 2003.

- [14] F. Rothlauf. *Representations for Genetic and Evolutionary Algorithms*. Physica-Verlag, 2002.
- [15] P. F. Stadler. Landscapes and their correlation functions. Technical Report 95-07-067, Santa Fe Institute, USA, 1995.
- [16] P. F. Stadler and G. P. Wagner. Algebraic theory of recombination spaces. *Evolutionary Computation*, 5(3):241–275, 1997.
- [17] A. M. Sutton. Personal Communication at Dagstuhl, 2010.
- [18] A. M. Sutton, L. D. Whitley, and A. E. Howe. A polynomial time computation of the exact correlation structure of k-satisfiability landscapes. In *GECCO '09: Proceedings of the 11th Annual conference on Genetic and evolutionary computation*, pages 365–372, Montreal, 2009. ACM.
- [19] V. K. Vassilev, T. C. Fogarty, and J. F. Miller. Smoothness, ruggedness and neutrality of fitness landscapes: from theory to application. In A. Ghosh and S. Tsutsui, editors, *Advances in evolutionary computing: theory and applications*, pages 3–44. Springer-Verlag New York, Inc., 2003.
- [20] D. Whitley, D. Hains, and A. Howe. Tunneling between optima: partition crossover for the traveling salesman problem. In *GECCO '09: Proceedings of the 11th Annual conference on Genetic and evolutionary computation*, pages 915–922, Montreal, 2009. ACM.
- [21] D. Whitley, A. M. Sutton, and A. E. Howe. Understanding elementary landscapes. In M. Keijzer, G. Antoniol, C. B. Congdon, K. Deb, B. Doerr, N. Hansen, J. H. Holmes, G. S. Hornby, D. Howard, J. Kennedy, S. Kumar, F. G. Lobo, J. F. Miller, J. Moore, F. Neumann, M. Pelikan, J. Pollack, K. Sastri, K. Stanley, A. Stoica, E.-G. Talbi, and I. Wegener, editors, *GECCO '08: Proceedings of the 10th annual conference on Genetic and evolutionary computation*, pages 585–592, Atlanta, GA, USA, 12-16 July 2008. ACM.
- [22] T. Yu and J. F. Miller. Through the interaction of neutral and adaptive mutations, evolutionary search finds a way. *Artificial Life*, 12(4):525–551, Fall 2006.