

# Instructions for using `Mutation_Age_estimation.R`

Luke C. Gandolfo\*, Melanie Bahlo† and Terence P. Speed‡

This document explains how to use the R script `Mutation_Age_estimation.R` to produce mutation age estimates and confidence intervals using the method described in [Gandolfo, Bahlo & Speed \(2014\)](#)

## 1 Obtain the data

Firstly, use the following procedure to obtain the *genetic length* of ancestral segments in sampled individuals (see [Figure 1](#)):

1. Line up your high density SNP data around the mutation locus.
2. Highlight *continuous* haplotype sharing between *two or more* individuals as you move away from the mutation locus. This sharing defines the ancestral segments.
3. Using a genetic map, record the genetic length of the segments. Specifically: record the segment “arm” lengths for each individual, i.e. the left arm length is the genetic distance between the mutation and the outermost marker heading in the direction of decreasing map position and the right arm length is the distance between the mutation and the outermost marker heading in the direction of increasing map position. Record the lengths in units of centiMorgans (cM).

To assess haplotype sharing for individuals with a single copy of the mutation (e.g. with a dominant mutation), or for individuals with copies of two distinct mutations at the same gene locus (i.e. with a compound heterozygous mutation), the data needs to be phased to resolve the underlying ancestral haplotype(s).<sup>1</sup> However, for individuals with two copies of the same recessive mutation, phasing is not required: ancestral segments are defined by continuous sharing of identical homozygous markers.<sup>2</sup>

---

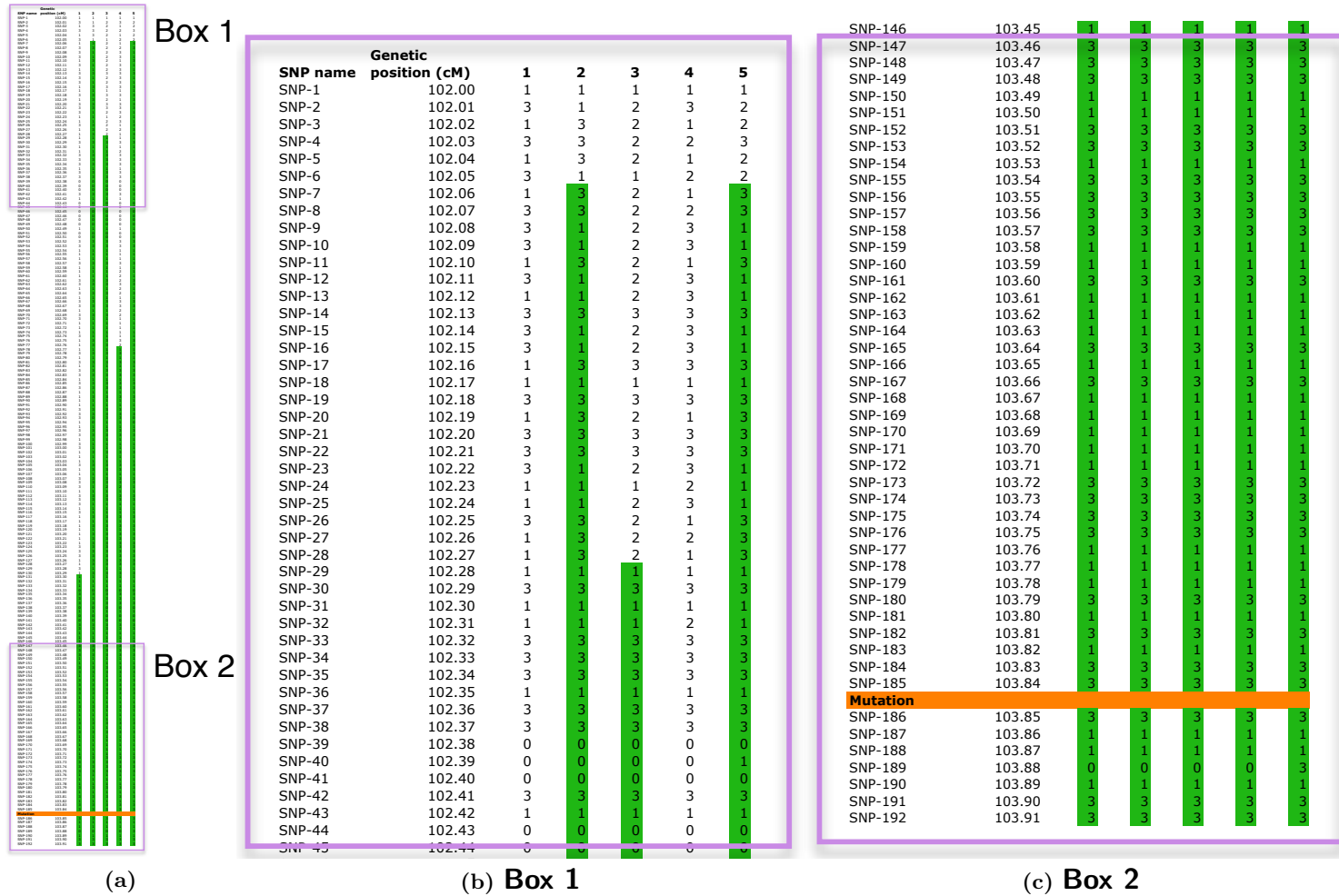
\*[gandolfo@wehi.edu.au](mailto:gandolfo@wehi.edu.au), or [gandolfo@student.unimelb.edu.au](mailto:gandolfo@student.unimelb.edu.au)

†[bahlo@wehi.edu.au](mailto:bahlo@wehi.edu.au)

‡[terry@wehi.edu.au](mailto:terry@wehi.edu.au)

<sup>1</sup>For an extensive discussion of phasing methods see the review by [Browning & Browning \(2011\)](#).

<sup>2</sup>To prepare the data in this case we recommend using the Perl script LINKDATAGEN ([Bahlo & Bromhead, 2009](#)), selecting the *complete data* option and then opening the resulting output in Excel.



**Figure 1:** An invented SNP data set on five individuals, lined up around the mutation (locus highlighted *orange*), with continuous haplotype sharing between two or more individuals highlighted *green*, which we infer to be ancestral segments: (a) shows segment arms to the left of the mutation (i.e. heading in the direction of decreasing map position); (b) is a magnified view of Box 1; and (c) is a magnified view of Box 2. Numbers code for genotypes: 1 = AA, 2 = AB, 3 = BB, and 0 = missing. Here we assume each individual has two copies of the same recessive mutation, thus the ancestral segments are defined by continuous sharing of identical homozygous markers. The left segment arm lengths are the genetic distance between the mutation and the outermost markers (the right arm lengths are found similarly).

Note that mutations and errors in marker data should be obvious: continuous sharing interrupted by one of these features will typically be followed by another large run of continuous sharing. When these features are obvious, they should be ignored. These features may not be obvious near segment ends, but in that case they have only a minor influence. Missing allele or genotype calls should also be ignored.

## 2 Run the script

Now take the following steps to execute the script in R:

1. Enter the length data:

```
> l.lengths = c(left arm length for individual 1, ..., left arm length
  for individual n)
> r.lengths = c(right arm length for individual 1, ..., right arm length
  for individual n)
```

Note that the components of each vector must correspond to the same individuals, i.e. the  $i^{\text{th}}$  position of each vector must refer to the same individual. Note also that the lengths must be in units of centiMorgans (cM).

2. Enter the desired parameters:

```
> confidence.coefficient = number between 0 and 1
> chance.sharing.correction = TRUE or FALSE
```

Confidence interval coverage is specified with `confidence.coefficient` (e.g. for 95% confidence intervals use 0.95). Set `chance.sharing.correction = TRUE` to correct for chance sharing, or `FALSE` otherwise. If this is set to `TRUE`, the following extra parameters are needed:

```
> median.allele.frequency = number between 0 and 1
> markers.on.chromosome = number
> length.of.chromosome = number
```

These parameters refer to the chromosome with the mutation: the first specifies the median population frequency of the A marker alleles (or the B marker alleles) across the chromosome; the second specifies the number of markers on the chromosome; and the third specifies the genetic length (in cM) of the chromosome.

3. Enter the following to execute the script:

```
> source(Mutation_Age_estimation.R)
```

Note that the script needs to be in the R working directory (either place the script there, or set the working directory to where the script is located).

### 3 Example

Entering the following length data (from a hypothetical sample of 6 individuals) and parameters:

```
> l.lengths = c(1.2, 3.4, 1.6, 4.8, 4.8, 2.3)
> r.lengths = c(2.0, 3.8, 2.7, 2.1, 3.8, 1.8)
> confidence.coefficient = 0.95
> chance.sharing.correction = TRUE
> median.allele.frequency = 0.7
> markers.on.chromosome = 68378
> length.of.chromosome = 290,
```

and running the script

```
> source(Mutation_Age_estimation.R),
```

produces the following output:

```
[1] Assuming an 'independent' genealogy: age estimate = 27.5 generations,
with confidence interval (15.5,49.2)
[1] Assuming a 'correlated' genealogy: age estimate = 22.3 generations,
with confidence interval (7.9,67.5)
```

### References

- Bahlo, M. & Bromhead, C. J. (2009). Generating linkage mapping files from Affymetrix SNP chip data. *Bioinformatics*, 25(15), 1961–1962.
- Browning, S. R. & Browning, B. L. (2011). Haplotype phasing: existing methods and new developments. *Nature Reviews Genetics*, 12(10), 703–714.
- Gandolfo, L. C., Bahlo, M., & Speed, T. P. (2014). Dating rare mutations from small samples with dense marker data. *Genetics*, 197(4), 1315–1327.