

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18

SimCo: a program to automate the comparison of multiple Structure runs

Clare A Rebbeck¹, Owen R. Jones^{1*}, Isheng J. Tsai¹

¹ Division of Biology, Imperial College London, Silwood Park Campus, Ascot, Berks. SL5 7PY, UK

* Corresponding author: Email: owen.jones@imperial.ac.uk

1 **ABSTRACT**

2 SUMMARY: The computer program "Structure" has become an increasingly popular tool
3 for scientists seeking to assign individuals to discrete populations using molecular data.
4 Typically, it is necessary to compare several runs of Structure outputs in order to evaluate the
5 consistency of assignments. A commonly used metric for comparing these assignments is the
6 similarity coefficient. Here we present a computer program, SimCo, which rigorously
7 calculates the similarity coefficient between multiple Structure runs. We provide versions that
8 can be implemented in PERL (SimCo-PERL) and R (SimCo-R) and make these publicly
9 available.

10 AVAILABILITY: The programs are available under the GNU General Public License and
11 are provided as supplementary information available from the journal's website. Simco for R is
12 available from <http://cran.r-project.org/>.

13 CONTACT: owen.jones@imperial.ac.uk

14 SUPPLEMENTARY INFORMATION: Available on Bioinformatics website

15

1 INTRODUCTION

2 Population geneticists often need to assign individuals to distinct populations. For
3 example, analysis on the genetic structure of domestic dog breeds (Parker et al. 2004) and
4 studies on human evolution (Adeyomo et al. 2005) and human migration (Falush et al. 2003)
5 focus on the evolutionary relationships of modern populations and the individuals within
6 populations. The program Structure (Pritchard et al. 2000; Falush et al. 2003) implements a
7 Bayesian model-based clustering algorithm which attempts to infer population structure and
8 assign individuals to populations probabilistically based on patterns of allele frequencies
9 [multilocus genotype data (e.g. SNPs, RFLPs and microsatellites etc.)]. A drawback of the
10 clustering method implemented in Structure is that it assumes K populations, where K is
11 defined by the user. If the value of K chosen is inappropriate, individuals can potentially form
12 weak, inconsistent associations with their assigned clusters, which can influence any inferences
13 drawn from an analysis (Pritchard et al. 2000). This, together with the inherent uncertainty in
14 the execution of Bayesian programs using the Markov Chain Monte Carlo simulation (Pearse
15 & Crandall 2004) make it advisable to perform multiple runs with the same data set.

16 Quantifying the similarity of outputs, by calculating the similarity coefficient is often
17 used as a measure of confidence for the populations (Parker et al. 2004, Rosenberg et al. 2002).
18 These calculations are, however, time consuming and complicated to carry out, especially with
19 increasing numbers of simulations, where the number of possible combinations increases
20 exponentially with the number of runs. To resolve this problem we have written programs,
21 which we call “SimCo”, implementable in R (R Development Core Team, 2005) and PERL, to
22 take the probabilistic assignments that Structure gives and carry out the necessary permutations
23 and similarity coefficient calculations automatically. We are releasing these programs under

1 the GNU General Public License. The programs are available as text files in the supplementary
2 information or from the authors.

3

4 **SYSTEM AND METHODS**

5 Structure output files contain data on the probability of assignment to one of K
6 populations. SimCo takes these output files and uses the assignment probabilities given therein
7 to calculate similarity coefficients as previously described by Rosenberg et al. (2002). It does
8 this for all possible pairwise combinations of n Structure runs ($n!/(2!(n-2)!)$) and then outputs
9 the summary statistics of the distribution of the calculated similarity coefficients (range,
10 moments, mean).

11 Because the population identities do not correspond with any particular output column
12 number, any two matrices cannot be directly compared. Therefore, the similarity coefficient
13 calculated within SimCo permutes the columns of the second matrix for every possible
14 combination of column ordering ($K!$ combinations), and calculates the absolute difference
15 between the first and second matrix for each permutation. The column-permutation with the
16 smallest difference is deemed to be the best fit for a direct comparison between the two
17 matrices.

18 **ALGORITHM**

19 The similarity coefficient for every pair of runs is calculated using equation 1;

$$20 \quad C(M1, M2) = 1 - \frac{\|M1 - M2\|_F}{\|M1 - 1/K\|_F}$$

21 Equation 1

1 where $C(M1, M2)$ is the similarity coefficient for the comparison between stochastic
2 matrices $M1$ and $M2$. The matrices have dimensions $I \times K$ (I = number of individuals; K =
3 number of supposed clusters that the individuals are assigned to) and are the outputs (i.e. the
4 probability of assignment to a particular cluster) from two separate Structure runs using the
5 same data. The notation $\|x\|_F$ is the Frobenius matrix norm (Golub and Van Loan 1996) and
6 $1/K$ is the $I \times K$ matrix with all the entries equal to $1/ K$.

7 **IMPLEMENTATION**

8 Both the Perl and R versions of the program can be used without the need to reformat
9 the Structure output file. Simply by editing the SimCo input command to include the file
10 pathways of simulations you wish to compare, any number of Structure runs can be analyzed
11 (memory dependent) simultaneously. Further details for running both the R and Perl versions
12 are given in PDF documents in the supplementary information. In addition there are help files
13 associated with the R version of simco (the ‘simco’ package). R is available for Mac, Windows
14 PC and Linux from <http://www.R-project.org>.

15

16 **DISCUSSION**

17 When implemented on a test data set with 3 runs, a K of 3 and 74 samples, both
18 SimCo-R and SimCo-PERL produced the same similarity coefficient value. Both programs
19 were fast with the test dataset completing the analysis in less than 1 second (on a Mac G5 with
20 4 x 2.5 GHz PowerPC processors and 2GB SDRAM). Tests indicate that, with datasets where a
21 larger number of runs and more assumed clusters are used, the program will take significantly
22 longer to run. We provide some test files as text documents in the supplementary information.

23

1 **CONCLUSION**

2 SimCo provides significant assistance to users of Structure, allowing the rigorous
3 comparison of the generated output files without the need to manipulate them. We are aware
4 that the need to carry out such matrix comparisons is widespread in biology (and the other
5 sciences). We are, therefore, making SimCo freely available on a license that allows free use,
6 distribution and alteration as long as we are acknowledged for our original work.

7

8 **ACKNOWLEDGEMENTS**

9 We thank Elaine Ostrander for providing us with a test dataset to work with. CAR and
10 IJT were supported by Wellcome Trust Studentships.

11 CONFLICT OF INTEREST: none declared

12

13 **REFERENCES**

- 14 Adeyomo, A. A. et al. (2005) Genetic structure in four West African population groups. *BMC*
15 *Genet.* 6:38, 1471-2156
- 16 Falush, D. et al. (2003) Traces of human migrations in *Helicobacter pylori* populations. *Science*.
17 299, 1582-1585.
- 18 Falush, D. et al. (2003) Inference of population structure using multilocus genotype data: linked
19 loci and correlated frequencies. *Genetics*. 164, 1567-1587.
- 20 Golub, G. H., and Van Loan, C. F. (1996) *Matrix computations*. Johns Hopkins University
21 Press, Baltimore.
- 22 Parker, G. H. et al. (2004) Genetic structure of the purebred domestic dog. *Science*. 304, 1160-
23 1164.

- 1 Pearse, D. E. & Crandall, K. A. (2004) Beyond FST: Analysis of population genetic data for
2 conservation. *Conserv. Genet.* 5, 585-602.
- 3 Pritchard, J. K. et al. (2000) Inference of population structure using multilocus genotype data.
4 *Genetics.* 155, 945-959.
- 5 R Development Core Team (2005) R: A language and environment for statistical computing. R
6 Foundation for Statistical Computing, Vienna, Austria. ISBN 3-900051-07-0, URL
7 <http://www.R-project.org>.
- 8 Rosenberg, N. A. et al. (2002) The genetic structure of human populations. *Science.* 298, 2381-
9 2385.