

Probabilistic graph models for landscape genetics

Brook G. Milligan¹

¹Department of Biology, New Mexico State University, Las Cruces, New Mexico 88003 USA

Corresponding author:

Brook G. Milligan¹

Email address: brook@nmsu.edu

ABSTRACT

Progress in landscape genetics depends on a strong conceptual foundation and the means of identifying mechanistic connections between environmental factors, landscape features, and genetic or genomic variation. Many existing approaches and much of the software commonly in use was developed for population genetics or statistics and is not entirely appropriate for landscape genetics. Probabilistic graph models provide a statistically rigorous and flexible means of constructing models directly applicable to landscape genetics. Sophisticated software also exists for the analysis of graph models. However, much of that software does not handle the types of data used for landscape genetics, model structures involving autoregressive spatial interaction between variables, or the scale of landscape genetics problems. Thus, an important priority for the field is to develop suitably flexible software tools for graph models that overcome these problems and allow landscape geneticists to explore meaningfully mechanistic and flexible models. We are developing such a library and applying it to examples in landscape genetics.

Keywords: landscape genetics, graph models, Bayesian inference, open source software

INTRODUCTION

One recurring theme in the landscape genetics literature is that progress is limited by available analytical methods (Balkenhol et al., 2009, 2016a; Guillot et al., 2009). In part this derives from the fact that many of the available analytical tools and much of the usable software were originally developed for population genetics or even broader statistical applications. They may include assumptions or be applicable to data that are not completely appropriate for landscape genetics studies. Because of this gap, there is no consensus in the literature regarding how to approach landscape genetics analysis (Balkenhol et al., 2016a). Nevertheless, development of a more comprehensive theory will come in part from an improved foundation of computational tools, especially open source ones, allowing explicit and flexible modeling.

This brief review focuses on three themes. First, it identifies the types of models most likely to contribute to advancement of a comprehensive theory of landscape genetics, improved mechanistic understanding, and better predictive power upon which, for example, conservation policy and management can be based. Second, it considers a set of open source software upon which such models for landscape genetics might be based. All of these turn out to have significant limitations.

Consequently, it also suggests characteristics that will be essential for the ongoing development of models and computational tools most likely to advance landscape genetics.

LANDSCAPE GENETICS AND BAYESIAN INFERENCE

The prevailing challenge in landscape genetics is identifying the mechanisms by which landscape and environmental factors influence genetic and genomic variation. More precisely, the central question in landscape genetics is the following: given data on intraspecific genetic variation across landscapes, what inferences are possible regarding the functional mechanisms and factors causing that variation? Framing the question in this way emphasizes the inherent connection between the science of landscape genetics and the nature of Bayesian inference.

BAYESIAN MODELS IN LANDSCAPE GENETICS

The natural connection between landscape genetics and Bayesian inference has led to the development of a variety of widely used Bayesian analysis methods. Although originally designed for population genetics, the most widely used is Structure, which identifies putative populations and assigns individuals to them (Pritchard et al., 2000). A second set seeks to identify population clusters by modeling allele frequency distributions in a spatially explicit way (Chen et al., 2007; Guillot et al., 2005a,b). More recently, Bayesian models that explicitly relate environmental gradients to spatially explicit allele frequency distributions have been developed (Coop et al., 2010; Frichot et al., 2013). One element is common to all of these models and associated software: each one covers a particular type of model and provides very limited opportunity for exploring related models or for expanding their scope. This is a serious limitation for a scientific field that repeatedly asserts that more mechanistic and predictive models and a stronger theoretical foundation is essential (Andrew et al., 2013; Balkenhol et al., 2016b; Guillot et al., 2009; Manel and Holderegger, 2013).

PROBABILISTIC GRAPH MODELS

This gap is not for lack of a general statistical framework that is completely applicable. Probabilistic graph models (Bishop, 2006; Koller and Friedman, 2009) are the means of describing and analyzing a broad range of models and sophisticated software exists to handle them. They are composed of random variables (vertices) and relationships between them (edges). Despite the superficial similarity involving graphs, probabilistic graph models are completely distinct from graph theory as applied to landscape genetics (Murphy et al., 2016). Generally, there is great scope for constructing general theories based upon manipulating probabilistic graph models to reflect interesting biological models within landscape genetics. However, software tools must exist that enable manipulation and analysis of the graphs, and the types of graphs available must match those required by landscape genetics. For many applications two types of graphs are enough: Bayesian networks represented by directed acyclic graphs (DAGs) and Markov random fields represented by undirected graphs. Landscape genetics models, however, often require more general types of graphs to accommodate, for example, spatially autoregressive relationships among random variables.

OPEN-SOURCE PROBABILISTIC GRAPH MODELS

While probabilistic graph models applied to landscape genetics do not generally harness their full flexibility, there exist modeling software that does. The most widely used is based upon the BUGS language for describing graph models, and includes WinBUGS, OpenBugs (Lunn et al., 2009) and JAGS (Plummer, 2015). The BUGS language allows textual description of general graph models that include a broad range of distributions. The textual description is translated into executable code, a process that introduces some of the limitations common to this type of modeling software. Another general graph modeling system is Stan, named for Stanislaw Ulam, an inventor of Monte Carlo approaches to inference (Carpenter et al., 2015; Gelman et al., 2015). Although more flexible in some ways than BUGS, Stan suffers from some of the same limitations that reduce its applicability to landscape genetics. It has the same limited data types and the execution environment is likewise limited by the Stan language. In addition to these two major classes of graph modeling software, a broad range of more specialized software systems is also available; many of these are summarized by Murphy (2014). Some are open source and may have potential for landscape genetics applications. However, they often handle a more limited range of graphs than is needed for landscape genetics, the data types are not well suited to landscape genetics, or their execution environments are limiting. This means that landscape geneticists face a fundamental challenge hindering development of a strong conceptual foundation for the field based upon the expressive power, flexibility, and statistical rigor of probabilistic graph models. Existing frameworks such as provided by BUGS, JAGS, and Stan offer much flexibility and power but are designed for types of graphs, random variables, and data types that are not ideally suited to landscape genetics. Other software libraries may suffer from these same limitations but in addition are much more difficult to program and well beyond the reach of typical landscape geneticists.

DESIGNING A PROBABILISTIC GRAPH MODEL FOR LANDSCAPE GENETICS

What then is the ideal design of a software system intended to harness the power, flexibility, and rigor of probabilistic graph models applied to landscape genetics? First and foremost, it must support a full range of relevant graph types, which in particular means not being limited to directed acyclic graphs. Second, it must support a full range of useful data types that landscape geneticists work with; in addition to simple scalars, vectors, and matrices, these include named alleles and genotypes, loci and chromosomes, spatial data of various sorts, and geographic locations. Ideally, user-defined or third-party data types should be easy to accommodate. Third, the algorithms available should be extensible to allow improved efficiency as needed. Fourth, the execution environment should not be limited to that encapsulated within a single, predefined program. This is especially important for landscape genetics models that may well encompass thousands or millions of random variables. Finally, all of this power and flexibility must be abstracted enough that a full spectrum of landscape geneticists can create simple models easily, test alternative and biologically relevant models flexibly, and improve upon the models and algorithms as needed. It is little surprise that existing software tools are unable to meet these stringent demands; they are largely conflicting and impossible to resolve without advanced software design. The most likely path forward (Lunn et al., 2009) leverages the power of C++ to present high-level abstractions based upon embedded domain specific languages

(de Guzman and Kaiser, 2016; Niebler, 2016) assembled with expression templates (Niebler, 2016; Veldhuizen, 1995) from highly reusable generic components (Stepanov and Rose, 2014). Although beyond the scope of this paper, we are following these design principles to implement a software library intended to provide the expressive power and computational performance demanded for advancing a coherent conceptual foundation for landscape genetics.

CONCLUSION

Landscape genetics suffers greatly from the absence of an analytical foundation that encourages development of a mechanistic understanding of the impact of environmental and landscape factors on genetic and genomic variation (Balkenhol et al., 2016a). This stems in part from the adoption of software tools and methods originally developed for other purposes. There exist well-established concepts and statistical approaches associated with probabilistic graph models that are ideally suited as the needed foundation for landscape genetics. Unfortunately, the associated software tools cannot be borrowed directly, because they are limited in ways that do not accommodate the needs of landscape geneticists. One priority that would directly advance the field and resolve these problems is the development of probabilistic graph model tools that do apply to landscape genetics. Despite the inherent difficulty of this task, we have developed a suitable library and are beginning to apply it to landscape genetics.

REFERENCES

- Andrew, R. L., Bernatchez, L., Bonin, A. L., Buerkle, C. A., Carstens, B. C., Emerson, B. C., Garant, D., Giraud, T., Kane, N. C., Rogers, S. M., Slate, J., Smith, H., Sork, V. L., Stone, G. N., Vines, T. H., Waits, L., Widmer, A., and Rieseberg, L. H. (2013). A road map for molecular ecology. *Molecular Ecology*, 22:2605–2626.
- Balkenhol, N., Cushman, S. A., Storfer, A. T., and Waits, L. P., editors (2016a). *Landscape genetics: concepts, methods, applications*. Wiley Blackwell, Hoboken, New Jersey.
- Balkenhol, N., Cushman, S. A., Waits, L. P., and Storfer, A. (2016b). Current status, future opportunities, and remaining challenges in landscape genetics. In Balkenhol, N., Cushman, S. A., Storfer, A. T., and Waits, L. P., editors, *Landscape genetics: concepts, methods, applications*, chapter 14, pages 247–255. Wiley Blackwell, Hoboken, New Jersey.
- Balkenhol, N., Gugerli, F., Cushman, S. A., Waits, L. P., Coulon, A., Arntzen, J. W., Holderegger, R., Wagner, H. H., and Participants of the Landscape Genetics Research Agenda Workshop 2007 (2009). Identifying future research needs in landscape genetics: where to from here? *Landscape Ecology*, 24:455–463.
- Bishop, C. M. (2006). *Pattern recognition and machine learning*. Springer, New York, New York.
- Carpenter, B., Lee, D., Gelman, A., Goodrich, B., Guo, J., Hoffman, M., Betancourt, M., Li, P., Brubaker, M. A., and Riddell, A. (2015). Stan: a probabilistic programming language. *Journal of Statistical Software*.
- Chen, C., Durand, E., Forbes, F., and François, O. (2007). Bayesian clustering algorithms ascertaining spatial population structure: a new computer program and a comparison study. *Molecular Ecology Notes*, 7:747–756.

- Coop, G., Witonsky, D., Di Rienzo, A., and Pritchard, J. K. (2010). Using environmental correlations to identify loci underlying local adaptation. *Genetics*, 185:1411–1423.
- de Guzman, J. and Kaiser, H. (2016). The Boost Spirit library. http://www.boost.org/doc/libs/1_61_0/libs/spirit/doc/html/index.html.
- Frichot, E., Schoville, S. D., Bouchard, G., and François, O. (2013). Testing for associations between loci and environmental gradients using latent factor mixed models. *Molecular Biology and Evolution*, 30:1687–1699.
- Gelman, A., Lee, D., and Guo, J. (2015). Stan: A probabilistic programming language for Bayesian inference and optimization.
- Guillot, G., Estoup, A., Mortier, F., and Cosson, J. F. (2005a). A spatial statistical model for landscape genetics. *Genetics*, 170:1261–1280.
- Guillot, G., Leblois, R., Coulon, A., and Frantz, A. C. (2009). Statistical methods in spatial genetics. *Molecular Ecology*, 18:4734–4756.
- Guillot, G., Mortier, F., and Estoup, A. (2005b). GENELAND: a computer package for landscape genetics. *Molecular Ecology Notes*, 5:712–715.
- Koller, D. and Friedman, N. (2009). *Probabilistic graphical models: principles and techniques*. MIT Press, Cambridge, MA.
- Lunn, D., Spiegelhalter, D., Thomas, A., and Best, N. (2009). The BUGS project: evolution, critique and future directions. *Statistics in Medicine*, 28:3049–3067.
- Manel, S. and Holderegger, R. (2013). Ten years of landscape genetics. *Trends in Ecology and Evolution*, 28:614–621.
- Murphy, K. (2014). Software packages for graphical models. <https://www.cs.ubc.ca/~murphyk/Software/bnsoft.html>.
- Murphy, M., Dyer, R., and Cushman, S. A. (2016). Graph theory and network models in landscape genetics. In Balkenhol, N., Cushman, S. A., Storfer, A. T., and Waits, L. P., editors, *Landscape genetics: concepts, methods, applications*, chapter 10, pages 165–179. Wiley Blackwell, Hoboken, New Jersey.
- Niebler, E. (2016). The Boost Proto library. http://www.boost.org/doc/libs/1_61_0/doc/html/proto.html.
- Plummer, M. (2015). JAGS version 4.0.0 user manual. Technical report.
- Pritchard, J. K., Stephens, M., and Donnelly, P. (2000). Inference of population structure using multilocus genotype data. *Genetics*, 155:945–959.
- Stepanov, A. A. and Rose, D. E. (2014). *From Mathematics to Generic Programming*. Addison-Wesley, first edition.
- Veldhuizen, T. (1995). Expression templates. *C++ Report*, 7(5):26–31.