

Probabilistic graph models for landscape genetics

Brook G. Milligan
Department of Biology
New Mexico State University
Las Cruces, New Mexico 88003 USA
brook@nmsu.edu

September 12, 2016

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.

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.

27 This brief review focuses on three themes. First, it identifies the types of models most likely to
28 contribute to advancement of a comprehensive theory of landscape genetics, improved mechanistic
29 understanding, and better predictive power upon which, for example, conservation policy and
30 management can be based. Second, it considers a set of open source software upon which such
31 models for landscape genetics might be based. All of these turn out to have significant limitations.
32 Consequently, it also suggests characteristics that will be essential for the ongoing development of
33 models and computational tools most likely to advance landscape genetics.

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

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

54 **Probabilistic graph models** This gap is not for lack of a general statistical framework that is
55 completely applicable. Probabilistic graph models (Bishop, 2006; Koller and Friedman, 2009) are
56 the means of describing and analyzing a broad range of models and sophisticated software ex-
57 ists to handle them. They are composed of random variables (vertices) and relationships between
58 them (edges). Despite the superficial similarity involving graphs, probabilistic graph models are
59 completely distinct from graph theory as applied to landscape genetics (Murphy et al., 2016). Gen-
60 erally, there is great scope for constructing general theories based upon manipulating probabilistic
61 graph models to reflect interesting biological models within landscape genetics. However, software
62 tools must exist that enable manipulation and analysis of the graphs, and the types of graphs avail-
63 able must match those required by landscape genetics. For many applications two types of graphs

64 are enough: Bayesian networks represented by directed acyclic graphs (DAGs) and Markov ran-
65 dom fields represented by undirected graphs. Landscape genetics models, however, often require
66 more general types of graphs to accommodate, for example, spatially autoregressive relationships
67 among random variables.

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

91 **Designing a probabilistic graph model for landscape genetics** What then is the ideal design of
92 a software system intended to harness the power, flexibility, and rigor of probabilistic graph models
93 applied to landscape genetics? First and foremost, it must support a full range of relevant graph
94 types, which in particular means not being limited to directed acyclic graphs. Second, it must
95 support a full range of useful data types that landscape geneticists work with; in addition to simple
96 scalars, vectors, and matrices, these include named alleles and genotypes, loci and chromosomes,
97 spatial data of various sorts, and geographic locations. Ideally, user-defined or third-party data
98 types should be easy to accommodate. Third, the algorithms available should be extensible to
99 allow improved efficiency as needed. Fourth, the execution environment should not be limited to
100 that encapsulated within a single, predefined program. This is especially important for landscape
101 genetics models that may well encompass thousands or millions of random variables. Finally, all of

102 this power and flexibility must be abstracted enough that a full spectrum of landscape geneticists
103 can create simple models easily, test alternative and biologically relevant models flexibly, and
104 improve upon the models and algorithms as needed. It is little surprise that existing software tools
105 are unable to meet these stringent demands; they are largely conflicting and impossible to resolve
106 without advanced software design. The most likely path forward (Lunn et al., 2009) leverages the
107 power of C++ to present high-level abstractions based upon embedded domain specific languages
108 (de Guzman and Kaiser, 2016; Niebler, 2016) assembled with expression templates (Niebler, 2016;
109 Veldhuizen, 1995) from highly reusable generic components (Stepanov and Rose, 2014). Although
110 beyond the scope of this paper, we are following these design principles to implement a software
111 library intended to provide the expressive power and computational performance demanded for
112 advancing a coherent conceptual foundation for landscape genetics.

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

124 References

- 125 Andrew, R. L., L. Bernatchez, A. L. Bonin, C. A. Buerkle, B. C. Carstens, B. C. Emerson,
126 D. Garant, T. Giraud, N. C. Kane, S. M. Rogers, J. Slate, H. Smith, V. L. Sork, G. N. Stone,
127 T. H. Vines, L. Waits, A. Widmer, and L. H. Rieseberg. 2013. A road map for molecular
128 ecology. *Molecular Ecology*, 22:2605–2626.
- 129 Balkenhol, N., S. A. Cushman, A. T. Storfer, and L. P. Waits, editors. 2016a. *Landscape genetics:
130 concepts, methods, applications*. Wiley Blackwell, Hoboken, New Jersey.
- 131 Balkenhol, N., S. A. Cushman, L. P. Waits, and A. Storfer. 2016b. Current status, future op-
132 portunities, and remaining challenges in landscape genetics. In Balkenhol, N., S. A. Cushman,
133 A. T. Storfer, and L. P. Waits, editors, *Landscape genetics: concepts, methods, applications*,
134 chapter 14, pages 247–255. Wiley Blackwell, Hoboken, New Jersey.
- 135 Balkenhol, N., F. Gugerli, S. A. Cushman, L. P. Waits, A. Coulon, J. W. Arntzen, R. Holderegger,
136 H. H. Wagner, and Participants of the Landscape Genetics Research Agenda Workshop 2007.

- 137 2009. Identifying future research needs in landscape genetics: where to from here? *Landscape*
138 *Ecology*, 24:455–463.
- 139 Bishop, C. M. 2006. *Pattern recognition and machine learning*. Springer, New York, New York.
- 140 Carpenter, B., D. Lee, A. Gelman, B. Goodrich, J. Guo, M. Hoffman, M. Betancourt, P. Li, M. A.
141 Brubaker, and A. Riddell. 2015. Stan: a probabilistic programming language. *Journal of*
142 *Statistical Software*.
- 143 Chen, C., E. Durand, F. Forbes, and O. François. 2007. Bayesian clustering algorithms ascertain-
144 ing spatial population structure: a new computer program and a comparison study. *Molecular*
145 *Ecology Notes*, 7:747–756.
- 146 Coop, G., D. Witonsky, A. Di Rienzo, and J. K. Pritchard. 2010. Using environmental correlations
147 to identify loci underlying local adaptation. *Genetics*, 185:1411–1423.
- 148 de Guzman, J. and H. Kaiser. 2016. The Boost Spirit library. [http://www.boost.org/
149 doc/libs/1_61_0/libs/spirit/doc/html/index.html](http://www.boost.org/doc/libs/1_61_0/libs/spirit/doc/html/index.html).
- 150 Frichot, E., S. D. Schoville, G. Bouchard, and O. François. 2013. Testing for associations between
151 loci and environmental gradients using latent factor mixed models. *Molecular Biology and*
152 *Evolution*, 30:1687–1699.
- 153 Gelman, A., D. Lee, and J. Guo. 2015. Stan: A probabilistic programming language for Bayesian
154 inference and optimization.
- 155 Guillot, G., A. Estoup, F. Mortier, and J. F. Cosson. 2005a. A spatial statistical model for landscape
156 genetics. *Genetics*, 170:1261–1280.
- 157 Guillot, G., R. Leblois, A. Coulon, and A. C. Frantz. 2009. Statistical methods in spatial genetics.
158 *Molecular Ecology*, 18:4734–4756.
- 159 Guillot, G., F. Mortier, and A. Estoup. 2005b. GENELAND: a computer package for landscape
160 genetics. *Molecular Ecology Notes*, 5:712–715.
- 161 Koller, D. and N. Friedman. 2009. *Probabilistic graphical models: principles and techniques*.
162 MIT Press, Cambridge, MA.
- 163 Lunn, D., D. Spiegelhalter, A. Thomas, and N. Best. 2009. The BUGS project: evolution, critique
164 and future directions. *Statistics in Medicine*, 28:30493067.
- 165 Manel, S. and R. Holderegger. 2013. Ten years of landscape genetics. *Trends in Ecology and*
166 *Evolution*, 28:614–621.
- 167 Murphy, K. June 16 2014. Software packages for graphical models. [https://www.cs.ubc.
168 ca/~murphyk/Software/bnsoft.html](https://www.cs.ubc.ca/~murphyk/Software/bnsoft.html).

- 169 Murphy, M., R. Dyer, and S. A. Cushman. 2016. Graph theory and network models in landscape
170 genetics. In Balkenhol, N., S. A. Cushman, A. T. Storfer, and L. P. Waits, editors, *Landscape ge-*
171 *netics: concepts, methods, applications*, chapter 10, pages 165–179. Wiley Blackwell, Hoboken,
172 New Jersey.
- 173 Niebler, E. 2016. The Boost Proto library. http://www.boost.org/doc/libs/1_61_0/doc/html/proto.html.
- 175 Plummer, M. October 1 2015. JAGS version 4.0.0 user manual. Technical report.
- 176 Pritchard, J. K., M. Stephens, and P. Donnelly. 2000. Inference of population structure using
177 multilocus genotype data. *Genetics*, 155:945–959.
- 178 Stepanov, A. A. and D. E. Rose. 2014. *From Mathematics to Generic Programming*. Addison-
179 Wesley, first edition.
- 180 Veldhuizen, T. June 1995. Expression templates. *C++ Report*, 7(5):26–31.