

ALGEBRAIC METHODS IN PHYLOGENETICS

MARTA CASANELLAS AND JOHN RHODES

To those outside the field, and even to some focused on empirical applications, phylogenetics may appear to have little to do with algebra. Probability and statistics are clearly important ingredients, as modeling and inferring evolutionary relationships motivate the field. Combinatorics is also an obvious component, as the graph-theoretic notions of trees, and more recently networks, are used to describe the relationships. But where does the algebra arise?

The models used in phylogenetics are necessarily complex. At the simplest they depend on a tree structure, as well as Markov matrices describing changes in nucleotide sequences along the edges. These two components result in probability distributions given by rather complicated polynomials on the parameters of the models, whose precise form reflects the structure of the tree. Even following standard statistical paradigms for inference, efficient calculation, such as by the Felsenstein pruning algorithm [Fel81] used in likelihood calculations, depends on understanding this algebraic structure.

But in the late 1980s the algebraic structure also suggested alternative inference frameworks to some researchers. These included the *phylogenetic invariants* of Cavender and Felsenstein [CF87], and of Lake [Lak87], and the Hadamard transform framework of Hendy and his colleagues [HP89, HPS94]. While this early explicitly algebraic work resulted in a number of interesting mathematical explorations, perhaps culminated in Evans and Speed's invariants work [ES93], it had little impact on practical inference as simulation studies seldom showed good performance [Hue95].

In the early 2000's, works of Allman and Rhodes [AR03] and of Sturmfels and Sullivant [SS05] revived interest in invariants. Interest in applying algebraic perspectives to statistical problems, especially in computational biology, was exemplified by the book of Pachter and Sturmfels [PS05], which helped draw new researchers to the field. Of course algebra in statistics has been present from the beginning, such as in Pearson's work [Pea94], but as theoretical and computational tools

of algebra have developed, they had remained largely outside of the inference toolbox.

In recent years, algebraic methods have been crucial to advances in the theory of phylogenetic inference (in particular, parameter identifiability of phylogenetic models [AR09, ADR18]) and in new methods of tree reconstruction [FSC16], [CK14] that are competitive with traditional frameworks. The tools that have been used draw from algebraic geometry, commutative algebra, computational algebra and algebraic statistics as well as group representation theory and algebraic combinatorics.

The works in this volume showcase the varied directions in which algebra is playing a role in current phylogenetic research.

Algebraic varieties underly the investigation of mixture models by Gross et al., as well as the study of maximum likelihood inference using recently developed numerical algebraic geometry tools by Kosta and Kubjas. Sumner and Woodhams focus more tightly on the modeling of sequence evolution, and the algebraic origin of nicely structured models.

A number of works move beyond simple evolution on a tree. The multispecies coalescent model, which describes the biological process by which gene trees may differ from species trees, is analyzed by Disanto and Rosenberg with tools of algebraic combinatorics. Long and Kubatko also consider this model, greatly weakening the assumptions necessary to justify the invariant-based SVDquartets method of species tree inference. Durden and Sullivan give an identifiability result for a k -mer based distance under the coalescent.

Moving from trees to networks, Kim et al. investigate the impact of admixture on phylogenetic distances and tree reconstruction. Considering both the coalescent and the hybridization, Baños mixes algebraic and combinatorial approaches to show the identifiability of many network features from gene tree data.

Two works highlight other algebraic tools. Terauds and Sumner apply representation theory to study improving distance estimates based on gene order through maximum likelihood. Yoshida et al. bring tropical geometry and algebra to bear on summarizing collections of trees, through a new form of principal component analysis.

Finally, Huber et al.'s work highlights the role of submodularity, a concept appearing widely in combinatorics and optimization, while Wicke and Fischer address open questions on the Shapely value of trees.

REFERENCES

- [ADR18] E. S. Allman, J. H. Degnan, and J. A. Rhodes. Split probabilities and species tree inference under the multispecies coalescent model. *Bull. Math. Biol.*, 80(1):64–103, Jan 2018.
- [AR03] E. S. Allman and J. A. Rhodes. Phylogenetic invariants of the general Markov model of sequence mutation. *Math. Biosci.*, 186:113–144, 2003.
- [AR09] E. S. Allman and J. A. Rhodes. The identifiability of covarion models in phylogenetics. *IEEE ACM Trans. Comput. Biol. Bioinformatics*, 6:76–88, 2009.
- [CF87] J. A. Cavender and J. Felsenstein. Invariants of phylogenies in a simple case with discrete states. *J. Class.*, 4:57–71, 1987.
- [CK14] J. Chifman and L. Kubatko. Quartet inference from snp data under the coalescent model. *Bioinformatics*, 30(23):3317–3324, 2014.
- [ES93] S. N. Evans and T. P. Speed. Invariants of some probability models used in phylogenetic inference. *Ann. Stat.*, 21(1):355–377, 1993.
- [Fel81] J. Felsenstein. Evolutionary trees from DNA sequences: a maximum likelihood approach. *J. Mol. Evol.*, 17:368–376, 1981.
- [FSC16] J. Fernández-Sánchez and M. Casanellas. Invariant versus classical quartet inference when evolution is heterogeneous across sites and lineages. *Syst. Biol.*, 65(2):280–291, 2016.
- [HP89] M. D. Hendy and D. Penny. A framework for the quantitative study of evolutionary trees. *Syst. Zool.*, 38:297–309, 1989.
- [HPS94] M. D. Hendy, D. Penny, and M. Steel. A discrete Fourier analysis for evolutionary trees. *Proc. Natl. Acad. Sci.*, 91:3339–3343, 1994.
- [Hue95] J. P. Huelsenbeck. Performance of phylogenetic methods in simulation. *Syst. Biol.*, 44:17–48, 1995.
- [Lak87] J. A. Lake. A rate-independent technique for analysis of nucleic acid sequences: evolutionary parsimony. *Mol. Biol. Evol.*, 4:167–191, 1987.
- [Pea94] K. Pearson. Contributions to the mathematical theory of evolution. *Philosophical Transactions of the Royal Society of London A: Mathematical, Physical and Engineering Sciences*, 185:71–110, 1894.
- [PS05] L. Pachter and B. Sturmfels, editors. *Algebraic Statistics for Computational Biology*. Cambridge University Press, New York, NY, USA, 2005.
- [SS05] B. Sturmfels and S. Sullivant. Toric ideals of phylogenetic invariants. *J. Comput. Biol.*, 12:204–228, 2005.

DPT. MATEMTIQUES, UNIVERSITAT POLITCNICA DE CATALUNYA

DPT. MATHEMATICS, UNIVERSITY OF ALASKA, FAIRBANKS