

# RESEARCH STATEMENT

MICHAEL SCHWEINBERGER

Since the pioneering work of R.A. Fisher, C.R. Rao, J. Neyman, and others, the bulk of statistical research has focused on attributes  $(X_i, Y_i)$  of individual population members  $i$  and scenarios in which  $n \geq 2$  independent observations  $(X_1, Y_1), \dots, (X_n, Y_n) \stackrel{\text{ind}}{\sim} \mathbb{P}$  from a probability law  $\mathbb{P}$  are available. In more recent times, a mounting body of evidence has revealed that the world of the twenty-first century is interconnected and interdependent, underscored by recent events that started out as local problems and turned into global crises (e.g., pandemics, political and military conflicts, economic and financial crises). More often than not, such events are unique and cannot be replicated, and the data at hand are discrete and dependent. Despite the fact that the interconnected world of the twenty-first century affects the welfare of billions of people around the world, **statistical learning with theoretical guarantees from discrete and dependent attributes  $(X_i, Y_i)$  and connections  $Z_{i,j}$  is an underresearched area.**

*My research seeks to bridge the gap between statistical theory and the social sciences and other fields by providing interpretable models for discrete and dependent data supported by statistical theory, with a view to studying non-causal or causal relationships among discrete and dependent attributes  $(X_i, Y_i)$  under network interference  $Z_{i,j}$ .*

## Selected research accomplishments

The following research accomplishments are ordered from more recent to less recent, with the exception of the last research accomplishment.

**Statistical learning from discrete and dependent attributes  $(X_i, Y_i)$  under network interference  $Z_{i,j}$ .** To learn how the interconnected world of the twenty-first century affects individual and collective outcomes of interest, data on attributes  $(X_i, Y_i)$  and connections  $Z_{i,j}$  are needed. More often than not, such data are discrete and dependent, and independent replications from the same source are unavailable. In scenarios with discrete and dependent data, it is natural to base statistical learning on interpretable models that possess conditional independence properties and admit exponential-family representations of conditional or joint distributions. Such models can be viewed as extensions of regression models for discrete and dependent attributes  $(X_i, Y_i)$  and connections  $Z_{i,j}$ . That said, some of the world's leading probabilists and statisticians have expressed concern about the probabilistic behavior of simplistic versions of exponential-family models for dependent connections  $Z_{i,j}$ , and questioned whether statistical learning with theoretical guarantees is possible based on a single observation of dependent connections  $Z_{i,j}$  [e.g., 15, 2, 8, 4, 21]. These concerns are important, because exponential-family models  $Z_{i,j}$  are widely used in practice and are related to generalized linear models in statistics, Markov random fields in machine learning, Ising models and Gibbs measures in physics, and Boltzmann machines in artificial intelligence (for which Geoffrey E. Hinton was awarded the Nobel prize 2024 in physics).

In a decade-long and continuing sequence of lead-authored publications starting in 2011 (e.g., JASA [26], JRSSB [34], Annals of Statistics [40], Bernoulli [29], Statistical Science [35], arXiv:2012.07167 [46], arXiv:2410.07555 [10]), I have taken steps to address these concerns. Among other things, I have demonstrated that the absence of desirable properties of the models considered by [15, 2, 26, 8, 4, 21] can be overcome by leveraging additional structure (observed or unobserved) [34, 40, 46, 33]. In addition, I have shown that models for discrete and dependent data with  $p \rightarrow \infty$

parameters can be learned from a single observation of discrete and dependent data, without sacrificing computational scalability and theoretical guarantees [46, 40]. By comparison, the small body of existing statistical theory for exponential-family models of discrete and dependent data assumes that the number of parameters  $p$  is fixed and makes other restrictive assumptions that limit the scope of the theoretical results to classic models in physics and related models in spatial statistics, e.g., Ising models with  $p = 1$  or  $p = 2$  parameters [e.g., 25, 5, 3, 13]. By contrast, my research focuses on large classes of models of discrete and dependent data with  $p \rightarrow \infty$  parameters, which help study how the interconnected and interdependent world of the twenty-first century affects individual and collective outcomes of interest and come with theoretical guarantees.

**I developed the first stochastic block models with dependent edges** [34, 29, 36, 1, 36, 11, 9]. Stochastic block models are widely used for learning from network data who is close to whom. Stochastic block models with dependent edges within communities, first introduced in my 2015 publication [34], combine the advantages of stochastic block models (capturing who is close to whom) and regression models for dependent connections (capturing local dependencies among connections). My research team has developed scalable computational-statistical methods [1, 48, 36] along with statistical theory [29]. The Japanese company Sansan Inc. applied these methods to a professional network with  $\sim 240,000$  members [6].

**I developed one of the first latent space models and the first statistical approach to hierarchical community detection** [38]. Latent space models are popular alternatives to stochastic block models for learning from network data who is close to whom [17, 41, 28]. The ultrametric latent space models I introduced in [38] have intrinsic hierarchical structure and can be used for hierarchical community detection. I published them one year after the Euclidean latent space models of Hoff et al. [16], seven years before the hyperbolic latent space models with intrinsic hierarchical structure of Krioukov et al. [20], and nineteen years before the hierarchical community detection method of Bickel and collaborators [22] [see, e.g., 41, 28].

**I made early contributions to the first widely used temporal network models and the first joint probability models of connections and outcomes** [e.g., 39, 42, 27, 23, 30, 44]. These models have been used in hundreds of publications in the social and health sciences for learning whether similar behavior among connected individuals (e.g., substance abuse among friends) is due to (a) the influence of friends, (b) the tendency to select similar others as friends, or (c) both. My contributions include likelihood-based inference [42], uncertainty quantification [39], statistical tests [27, 23], latent variable models [30], and statistical software [44].

**To gain insight into the interconnected and interdependent world of the twenty-first century, I have designed stochastic models for real-world phenomena**, e.g., hate speech on social media [10], mental health [19], substance abuse [43], epidemics [32], neuroscience [35], air pollution [31], disaster response [37], terrorist networks [34], systemic risk in software networks [9], online trust networks [48], online educational assessments [19, 18], product recommendation [1], financial networks [27, 30], social networks [47], socio-economic segregation [24], political hierarchy [12], and sport analytics [14].

## Selected directions of future research

**Causal inference under interference.** At the heart of science is the question of cause and effect. I am interested in causal inference for attributes  $(X_i, Y_i)$  under interference  $Z_{i,j}$ . Interference arises when the outcomes of units are affected by the treatments or outcomes of other units. The resulting phenomenon is known as spillover: Treating a subset of units can affect the outcomes of

other units, including untreated and treated units. Understanding spillover is imperative in real-world applications: e.g., financial incentives can result in spillover that affecting market outcomes. Two forms of spillover can be distinguished: treatment spillover (the treatments  $X_j$  of other units  $j$  can affect the outcome  $Y_i$  of unit  $i$ ) and outcome spillover (the outcomes  $Y_j$  of other units  $j$  can affect the outcome  $Y_i$  of unit  $i$ ). The bulk of research has focused on treatment spillover, which implies that outcomes  $Y_1, \dots, Y_n$  are *independent conditional on treatments*  $X_1, \dots, X_n$ . My research team is interested in both treatment and outcome spillover, which implies that outcomes  $Y_1, \dots, Y_n$  are *dependent conditional on treatments*  $X_1, \dots, X_n$ . We seek to answer the following two questions:

(a) **Opening black boxes.** How can the indirect causal effect of treatments on outcomes due to treatment and outcome spillover be characterized as an explicit mathematical function of the direct treatment effect, the effect of treatment spillover, and the effect of outcome spillover? To date, the indirect causal effect due to treatment and outcome spillover resembles a black box, in the sense that it is unknown how it can be characterized as an explicit mathematical function of the treatment effect and the effect of treatment spillover and outcome spillover.

(b) **External validity.** How can conclusions based on a sample of outcomes be generalized to the population of interest, when the outcomes are dependent due to treatment and outcome spillover and therefore the observed outcomes depend on unobserved outcomes? While non-Bayesian or Bayesian data-augmentation approaches could be used, such approaches are time-consuming. Scalable approaches with theoretical guarantees are unknown.

The advances of my research team enable answers to these questions with theoretical guarantees.

**Any question about statistical procedures of attribute data  $(X_i, Y_i)$  can be asked about discrete and dependent attributes  $(X_i, Y_i)$  and connections  $Z_{i,j}$ .** Many of these questions are unanswered. My research team intends to answer them.

(a) **Model selection.** As a case in point, there are countless models of discrete and dependent data, but model selection procedures are scarce and lack either computational scalability or theoretical guarantees or both. My research team plans to work on a scalable approach to model selection in dependent-data problems with intractable likelihood functions. We intend to explore two directions of research, one based on pseudo-likelihood Dantzig selectors and the other one based on pseudo-likelihood Bayesian procedures, and provide theoretical guarantees.

(b) **Uncertainty quantification.** In applications, it is important to provide a disclaimer, acknowledging that statistical conclusions based on data are subject to error. In scenarios when the number of parameters is unbounded and a single observation of discrete and dependent random variables is available, it is not obvious how to quantify uncertainty, because the small- and large-sample distributions of many statistical quantities are unknown. To place uncertainty quantification and statistical tests on firm mathematical grounds, Berry-Esseen-type bounds for bounding the error of normality approximations for discrete and dependent data are needed. That said, few Berry-Esseen-type bounds for discrete and dependent data exist. The few existing Berry-Esseen-type bounds impose strong restrictions on dependence, such as local dependence [45, Theorem 2.5] or strong mixing conditions [7, Theorem 3.27, p. 34]. These restrictions may or may not be satisfied in social science applications. My research team intends to develop Berry-Esseen-type bounds under weaker restrictions on dependence.

**Stochastic models of network-space-time data.** Many real-world phenomena involve networks, space, and time. I intend to help data scientists design stochastic processes involving networks, space, and time that do justice to the complexity of the interconnected and interdependent world of the twenty-first century, expanding my work on stochastic models of network-space data and network-time data to network-space-time data.

## References

- [1] Babkin, S., Stewart, J. R., Long, X., and Schweinberger, M. (2020), “Large-scale estimation of random graph models with local dependence,” *Computational Statistics & Data Analysis*, 152, 1–19.
- [2] Bhamidi, S., Bresler, G., and Sly, A. (2011), “Mixing time of exponential random graphs,” *The Annals of Applied Probability*, 21, 2146–2170.
- [3] Chatterjee, S. (2007), “Estimation in spin glasses: A first step,” *The Annals of Statistics*, 35, 1931–1946.
- [4] Chatterjee, S., and Diaconis, P. (2013), “Estimating and understanding exponential random graph models,” *The Annals of Statistics*, 41, 2428–2461.
- [5] Comets, F. (1992), “On consistency of a class of estimators for exponential families of Markov random fields on the lattice,” *The Annals of Statistics*, 20, 455–468.
- [6] Dahbura, J. N. M., Komatsu, S., Nishida, T., and Mele, A. (2021), “A structural model of business card exchange networks,” *arXiv:2105.12704*, 1–33.
- [7] Dedecker, J., Doukhan, P., Lang, G., Leon, J. R., Louhichi, S., and Prieur, C. (eds.) (2007), *Weak Dependence: With Examples and Applications*, Springer-Verlag.
- [8] Fienberg, S. E. (2012), “A brief history of statistical models for network analysis and open challenges,” *Journal of Computational and Graphical Statistics*, 21, 825–839.
- [9] Fritz, C., Georg, C.-P., Mele, A., and Schweinberger, M. (2024), “A Strategic Model of Software Dependency Networks,” in *25th ACM (Association for Computing Machinery) Conference on Economics and Computation (EC ’24)*.
- [10] Fritz, C., Schweinberger, M., Bhadra, S., and Hunter, D. R. (2024), “A regression framework for studying relationships among attributes under network interference,” *arXiv:2410.07555*.
- [11] Fritz, C., Schweinberger, M., Komatsu, S., Martínez Dahbura, J. N., Nishida, T., and Mele, A. (2024), *bigergm: Fit, Simulate, and Diagnose Hierarchical Exponential-Family Models for Big Networks*, R package version 1.2.1.
- [12] Fritz, C., Yuan, Y., and Schweinberger, M. (2024), “Hyperbolic latent space models for hypergraphs,” in *Proceedings of the Twenty-Eighth International Conference on Artificial Intelligence and Statistics (AISTATS)*, submitted.
- [13] Ghosal, P., and Mukherjee, S. (2020), “Joint estimation of parameters in Ising model,” *The Annals of Statistics*, 48, 785–810.
- [14] Grieshop, N., Feng, Y., Hu, G., and Schweinberger, M. (2024), “A continuous-time stochastic process for high-resolution network data in sports,” *Statistica Sinica*, to appear.
- [15] Handcock, M. S. (2003), “Statistical Models for Social Networks: Inference and Degeneracy,” in *Dynamic Social Network Modeling and Analysis: Workshop Summary and Papers*, eds. Breiger, R., Carley, K., and Pattison, P., Washington, D.C.: National Academies Press, pp. 1–12.

- [16] Hoff, P. D., Raftery, A. E., and Handcock, M. S. (2002), “Latent space approaches to social network analysis,” *Journal of the American Statistical Association*, 97, 1090–1098.
- [17] Hunter, D. R., Krivitsky, P. N., and Schweinberger, M. (2012), “Computational statistical methods for social network models,” *Journal of Computational and Graphical Statistics*, 21, 856–882.
- [18] Jeon, M., Jin, I. H., Schweinberger, M., and Baugh, S. (2021), “Mapping unobserved item–respondent interactions: A latent space item response model with interaction map,” *Psychometrika*, 86, 378–403.
- [19] Jeon, M., and Schweinberger, M. (2024), “A latent process model for monitoring progress towards hard-to-measure targets, with applications to mental health and online educational assessments,” *The Annals of Applied Statistics*, 18, 2123–2146.
- [20] Krioukov, D., Papadopoulos, F., Kitsak, M., Vahdat, A., and Boguna, M. (2010), “Hyperbolic geometry of complex networks,” *Physical Review E*, 82.
- [21] Lauritzen, S., Rinaldo, A., and Sadeghi, K. (2018), “Random networks, graphical models and exchangeability,” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 80, 481–508.
- [22] Li, T., Lei, L., Bhattacharyya, S., van den Berge, K., Sarkar, P., Bickel, P. J., and Levina, E. (2022), “Hierarchical community detection by recursive partitioning,” *Journal of the American Statistical Association*, 117, 951–968.
- [23] Lospinoso, J., Schweinberger, M., Snijders, T. A. B., and Ripley, R. (2011), “Assessing and accounting for time heterogeneity in stochastic actor oriented models,” *Advances in Data Analysis and Classification*, 5, 147–176.
- [24] Nandy, S., Holan, S. H., and Schweinberger, M. (2024), “A socio-demographic latent space approach to spatial data when geography is important but not all-important,” *arXiv:2304.03331*.
- [25] Pickard, D. K. (1987), “Inference for discrete Markov fields: The simplest non-trivial case,” *Journal of the American Statistical Association*, 82, 90–96.
- [26] Schweinberger, M. (2011), “Instability, sensitivity, and degeneracy of discrete exponential families,” *Journal of the American Statistical Association*, 106, 1361–1370.
- [27] — (2012), “Statistical modeling of digraph panel data: goodness-of-fit,” *British Journal of Mathematical and Statistical Psychology*, 65, 263–281.
- [28] — (2019), “Random graphs,” in *Wiley StatsRef: Statistics Reference Online*, eds. Everitt, B., Molenberghs, G., Piegorsch, W., Ruggeri, F., Davidian, M., and Kenett, R., Wiley, pp. 1–11.
- [29] — (2020), “Consistent structure estimation of exponential-family random graph models with block structure,” *Bernoulli*, 26, 1205–1233.
- [30] — (2020), “Statistical inference for continuous-time Markov processes with block structure based on discrete-time network data,” *Statistica Neerlandica*, 74, 342–362.

- [31] Schweinberger, M., Babkin, S., and Ensor, K. B. (2017), “High-dimensional multivariate time series with additional structure,” *Journal of Computational and Graphical Statistics*, 26, 610–622.
- [32] Schweinberger, M., Bomiriya, R. P., and Babkin, S. (2022), “A semiparametric Bayesian approach to epidemics, with application to the spread of the coronavirus MERS in South Korea in 2015,” *Journal of Nonparametric Statistics*, 34, 628–662.
- [33] Schweinberger, M., and Fritz, C. (2023), “Invited discussion of “A tale of two datasets: Representativeness and generalisability of inference for samples of networks” by P.N. Krivitsky, P. Coletti, and N. Hens,” *Journal of the American Statistical Association*, 118, 2225–2227.
- [34] Schweinberger, M., and Handcock, M. S. (2015), “Local dependence in random graph models: characterization, properties and statistical inference,” *Journal of the Royal Statistical Society, Series B (Statistical Methodology)*, 77, 647–676.
- [35] Schweinberger, M., Krivitsky, P. N., Butts, C. T., and Stewart, J. R. (2020), “Exponential-family models of random graphs: Inference in finite, super, and infinite population scenarios,” *Statistical Science*, 35, 627–662.
- [36] Schweinberger, M., and Luna, P. (2018), “HERGM: Hierarchical exponential-family random graph models,” *Journal of Statistical Software*, 85, 1–39.
- [37] Schweinberger, M., Petrescu-Prahova, M., and Vu, D. Q. (2014), “Disaster response on September 11, 2001 through the lens of statistical network analysis,” *Social Networks*, 37, 42–55.
- [38] Schweinberger, M., and Snijders, T. A. B. (2003), “Settings in social networks: A measurement model,” *Sociological Methodology*, 33, 307–341.
- [39] Schweinberger, M., and Snijders, T. A. B. (2007), “Markov models for digraph panel data: Monte Carlo-based derivative estimation,” *Computational Statistics and Data Analysis*, 51, 4465–4483.
- [40] Schweinberger, M., and Stewart, J. R. (2020), “Concentration and consistency results for canonical and curved exponential-family models of random graphs,” *The Annals of Statistics*, 48, 374–396.
- [41] Smith, A. L., Asta, D. M., and Calder, C. A. (2019), “The geometry of continuous latent space models for network data,” *Statistical Science*, 34, 428–453.
- [42] Snijders, T. A. B., Koskinen, J., and Schweinberger, M. (2010), “Maximum likelihood estimation for social network dynamics,” *The Annals of Applied Statistics*, 4, 567–588.
- [43] Snijders, T. A. B., Steglich, C. E. G., and Schweinberger, M. (2007), “Modeling the co-evolution of networks and behavior,” in *Longitudinal models in the behavioral and related sciences*, eds. van Montfort, K., Oud, H., and Satorra, A., Lawrence Erlbaum, pp. 41–71.
- [44] Snijders, T. A. B., Steglich, C. E. G., Schweinberger, M., and Huisman, M. (2010), *Manual for Siena 3.0*, Department of Statistics, University of Oxford, UK.
- [45] Stewart, J. R. (2024), “Rates of convergence and normal approximations for estimators of local dependence random graph models,” *arXiv:2404.11464*.

- [46] Stewart, J. R., and Schweinberger, M. (2024), “Pseudo-likelihood-based  $M$ -estimators for random graphs with dependent edges and parameter vectors of increasing dimension,” *arXiv:2012.07167*.
- [47] Stewart, J. R., Schweinberger, M., Bojanowski, M., and Morris, M. (2019), “Multilevel network data facilitate statistical inference for curved ERGMs with geometrically weighted terms,” *Social Networks*, 59, 98–119.
- [48] Vu, D. Q., Hunter, D. R., and Schweinberger, M. (2013), “Model-based clustering of large networks,” *The Annals of Applied Statistics*, 7, 1010–1039.