

# RESEARCH STATEMENT

MICHAEL SCHWEINBERGER

My research is concerned with methodological, computational, and theoretical aspects of learning from

- discrete and dependent data without independent replications: e.g., network, spatial, and temporal data;
- structured data, that is, data with additional structure, either observed or unobserved: e.g., block structure, multilevel, spatial, and temporal structure;
- social science data: e.g., educational data, epidemiological data, and social network data.

My research has been facilitated by NWO award Rubicon-44606029 (sole PI), NSF award DMS-1513644 (sole PI), NSF award DMS-1812119 (sole PI), and ARO award W911NF-21-1-0335 (lead PI).

## Overview

My research is motivated by discrete and dependent data without independent replications, such as network, spatial, and temporal data. My ideas of how to learn from dependent data without independent replications are elaborated in one of the simplest possible settings: statistical exponential families (Wainwright and Jordan, Foundations and Trends in Machine Learning, 2008). Exponential families are widely used throughout data science, as stand-alone models or building blocks of more complex models. The fundamental role of exponential families in data science is exemplified by the prominent role of multivariate Gaussians, but there are numerous other examples, including generalized linear models, undirected graphical models with exponential-family parameterizations, random graph models with exponential-family parameterizations, Markov random fields in machine learning, and Boltzmann machines in artificial intelligence. For example, in statistical network analysis, random graph models with exponential-family parameterizations are widely used for specifying models that capture dependencies in network data. These statistical exponential-family models represent an alternative to latent variable models that induce dependence through latent variables, but do not provide data scientists with a simple and flexible approach to specifying and comparing models with competing network dependencies and network features of interest, despite the fact that any node-exchangeable random graph model can be expressed as a latent variable model by virtue of the Aldous-Hoover theorem (Bickel and Chen, PNAS, 2009).

## Selected highlight

Consider network data, which are dependent data without independent replications. Since the 1950s, social scientists have pointed out that connections depend on other connections: e.g., the frequent observation that “a friend of a friend is a friend” suggests that friendships are dependent. In applications, population probability models are learned from a single observation of a population network or subnetworks sampled from a population network. That raises an important question:

*What can we learn about an interconnected and interdependent world where connections depend on other connections, without having the benefit of independent observations from the same source?*

In a decade-long sequence of first- and single-authored papers (e.g., Annals of Statistics, 2020; Bernoulli, 2020; Statistical Science, 2020; Journal of the Royal Statistical Society, Series B, 2015; Journal of the American Statistical Association, 2011) and more recent papers (e.g., Stewart and Schweinberger, 2021), I have:

1. Studied the properties of ill-posed statistical exponential-family models of network data (e.g., model near-degeneracy). My work (Schweinberger, JASA, 2011) preceeded Chatterjee and Diaconis (AOS, 2013).
2. Shown how well-posed models of network data can be constructed based on statistical exponential families, with desirable properties.
3. Demonstrated that statistical learning of an unbounded number of parameters based on a single observation of dependent data from a statistical exponential family is possible, with theoretical guarantees.

4. Developed scalable methods for statistical learning of an unbounded number of parameters based on a single observation of dependent data from a statistical exponential family, with theoretical guarantees.

There is a common thread that connects these advances: **the importance of additional structure**. Models that lack mathematical structure to control the dependence among connections can be ill-posed, but endowing models with additional structure can help control dependence and result in well-posed models with desirable properties. In addition, weak dependence facilitates concentration-of-measure results, which in turn facilitate consistency results. In other words, endowing models with additional structure has two advantages:

1. It facilitates the construction of well-posed models with desirable properties.
2. It facilitates statistical learning with theoretical guarantees.

Examples of additional structure are block structure, multilevel, spatial, and temporal structure. Moreover, additional structure helps answer fundamental questions about the statistical analysis of dependent network data raised by leading probabilists (e.g., Chatterjee and Diaconis, AOS, 2013) and statisticians (e.g., Fienberg, JCGS, 2012), as discussed in Schweinberger et al. (Statistical Science, 2020).

## Selected directions of future research

**Online educational assessment data:** In collaboration with Minjeong Jeon (Graduate School of Education & Information Studies, University of California, Los Angeles), I am working on educational assessment data. Among other things, we are developing statistical interaction and learning progression maps, with a view to providing educators with visual tools for monitoring student progress and detecting (underrepresented) groups of students who need more, and different support than other students.

**Stochastic processes involving networks, space, and time:** Many real-world processes involve networks, space, and time: e.g., infectious diseases spread through contacts, contacts depend on geographical space, and contacts change over time. While there are existing stochastic processes indexed by networks, space, time or combinations of them, many of them make either simplifying assumptions or have unknown probabilistic and statistical properties. One of my directions of future research is to design stochastic processes indexed by networks, space, and time that do justice to the complexity of network-mediated phenomena, and develop scalable statistical and computational methods and theoretical guarantees for learning them from data.

**Scalable selection of models of dependent data without independent replications and intractable likelihood functions:** Developing scalable model selection procedures with theoretical guarantees is non-trivial when the likelihood function is intractable, the number of parameters is large, and the data consists of a single observation of dependent random variables. Such scenarios arise in the statistical analysis of discrete and dependent data, such as discrete network, spatial, and temporal data. For example, there are many models of dependent network data, but no scalable model selection procedures with theoretical guarantees are known. I am working on a scalable approach to model selection in dependent data problems with intractable likelihood functions based on pseudo- and composite-likelihood-based regularization methods.

**Quantifying uncertainty of statistical learning based on dependent data without independent replications:** In applications of data science, it is important to provide a disclaimer, acknowledging that statistical conclusions based on data are subject to error and quantifying the uncertainty about the conclusions. In scenarios when the number of parameters is unbounded and a single observation of discrete and dependent data is available, it is not obvious how to quantify uncertainty, because the distributions of many statistical quantities are unknown. A natural approach to capturing uncertainty is a Bayesian approach. While Bayesian approaches to dependent network data and other discrete and dependent data without independent replications exist, most of them are either not scalable or theoretical guarantees are unknown. I intend to elaborate on scalable Bayesian approaches to uncertainty quantification for discrete and dependent data without independent replications based on pseudo- and composite-likelihood functions, with theoretical guarantees.