

General measures of Attribution Disclosure Risk for gauging privacy of synthetic data

Yongjae Kim

Seoul National University

As the demand for synthetic data continues to grow, there is an increasing need for rigorous measures to assess whether synthetic data is safe or poses significant privacy risks. Correct Attribution Probability (CAP) is a widely used risk measure; however, its theoretical foundation has not been fully established within a solid statistical framework. In this paper, we propose a statistical framework for defining CAP and introduce a modified version to clarify its theoretical meaning. We also demonstrate the limitations of CAP as a comprehensive risk measure and argue why it cannot serve as an all-encompassing solution. Furthermore, we develop a generalized version of CAP, termed Attribution Disclosure Risk (ADR), which provides a more comprehensive and versatile assessment of synthetic data risk, incorporating CAP as a special case at both the population and sample levels. Numerical studies demonstrate that our proposed measure consistently captures the risk inherent in synthetic data and offers flexibility to accommodate various intruder scenarios, applicable to both simulated and real datasets.