

Differential privacy based on importance weighting

Zhanglong Ji · Charles Elkan

Received: 17 February 2013 / Accepted: 11 June 2013 / Published online: 28 June 2013
© The Author(s) 2013

Abstract This paper analyzes a novel method for publishing data while still protecting privacy. The method is based on computing weights that make an existing dataset, for which there are no confidentiality issues, analogous to the dataset that must be kept private. The existing dataset may be genuine but public already, or it may be synthetic. The weights are importance sampling weights, but to protect privacy, they are regularized and have noise added. The weights allow statistical queries to be answered approximately while provably guaranteeing differential privacy. We derive an expression for the asymptotic variance of the approximate answers. Experiments show that the new mechanism performs well even when the privacy budget is small, and when the public and private datasets are drawn from different populations.

Keywords Privacy · Differential privacy · Importance weighting

1 Introduction

Suppose that a hospital possesses a dataset concerning patients, their diseases, their treatments, and the outcomes of treatments. The hospital faces a fundamental conflict. On the one hand, to protect the privacy of the patients, the hospital wants to keep the dataset secret. On the other hand, to allow science to progress, the hospital wants to make the dataset public. This conflict is the issue addressed by research on privacy-preserving data mining. How can a data owner simultaneously both publish a dataset and conceal it?

We analyze here a new approach to resolving the fundamental tension between publishing and concealing data. The new approach is based on a mathematical technique called importance weighting that has proved to be valuable in several other areas of research (Hastings 1970). The essential idea is as follows. Let D be the set of records that the owner must

Editors: Hendrik Blockeel, Kristian Kersting, Siegfried Nijssen, and Filip Železný.

Z. Ji · C. Elkan (✉)

Department of Computer Science and Engineering 0404, University of California, San Diego, USA
e-mail: elkan@ucsd.edu