



Balancing privacy and utility when releasing genetic data



Carlos Rojo¹, Daniel Aranki², Ruzena Bajcsy²

1. San José City College, 2. UC Berkeley

Abstract

Genetic datasets are a rich information source that continue to guide biomedical research. However, genetic datasets are extremely sensitive, as they can be used to reveal information like ethnicity and disease risk. Here we explore the implementation of user-specified personalized privacy to balance the need to release genetic data with high utility (able to predict disease risk) with the desire to keep certain information private (such as ancestry).

0

Motivation and Goal

Single variations (SNPs) reveal disease risk and ancestry.

Goal: user-defined personal privacy on genetic data

Diabetes Classifier

Previous work has identified SNPs that are associated with Type 1 diabetes. A model to predict T1D based on genetic data established by Winkler, C., Krumsiek, J. et al. (1) Can use this weighted model to classify Type 1 diabetes risk based on genetic data of individual.



Increasingly being used in the clinic.

- Recent privacy breaches show that even releasing summary statistics is not private.
- How can we share genetic data while maintaining privacy?

"Which of the following would you like to hide when releasing your data?"

Alzheimer's Risk

Type I Diabetes Risk

Ancestry

$$\log_e(\frac{p}{1-p}) = \beta_o + \beta_1(SNP_1) + \beta_2(SNP_2) + \dots + \beta_n(SNP_n)$$

p =probability of having Type 1 diabetes B_0 = baseline Type 1 diabetes risk B_i = proportion that each SNP contributes to Type 1 diabetes risk SNP_{n} = Sequence at the SNP. That is, whether indiv. has disease allele

Analysis Overview



Results



Ancestry Classifier

- Principle Components Analysis (PCA) highlights principle SNPs that explain variability within genetic dataset.
- PCA can be used to stratify individuals based on ancestry.







Left: Figure demonstrating how ability to classify ethnicity changes when you remove ancestrydetermining SNPs (center colored clusters). Harder to distinguish ancestry via Euclidean distance (shown is what happens when you mask all SNPs except those that cause T1D and chromosome 6). **<u>Right</u>**: Accuracy of ethnicity classifier after removing all SNPs except those that determine T1D risk. Utility has not changed (since no T1D SNPs were removed) but privacy is increased.

Conclusions and Future Directions

- Whenever privatizing data, there can be a trade off between utility and privacy.
- Using genetic classifiers, personalized privacy on genetic data is possible. Here we demonstrate that you can effectively remove the prediction of ancestry while maintaining ability to diagnose a disease like Type 1 diabetes.
- Future work will attempt to increase the utility by allowing for prediction of risk for multiple diseases, while still privatizing ancestry determination.
- Future work will establish global guarantee against future classifiers.
- This work will also require addressing the issue of *linkage disequilibrium* (LD, or correlation between SNPs) to ensure privacy:



Removing SNP₂ from the dataset provides no privacy as it can be deduced from the sequence at SNP₁ and SNP₃. SNP₁, SNP₂, and SNP₃ are said to be in LD

The authors would like to thank the E3S center and its staff, Josephine Yuen, James Hake, and Lea Marlor for coordinating the RET summer program.



