# Differential Privacy and Synthetic Data Generation using PrivBayes

Jan Reiter Sørensen

# Agenda

- Problem statement

- Differential privacy

- The Laplace mechanism and the exponential mechanism

- PrivBayes

- Investigating synthetic data by PrivBayes

    - Generating some data

    - The mixture parameter

    - Overall privacy and power of Pearson's test

- An example

# Problem Statement boiled down

1. How are differentially private methods for synthesizing data defined and explained using a consistent and precise mathematical language?

2. How are differentially private synthesizers such as the smoothed histogram mechanism and PrivBayes implemented?
   a. In practice
   b. For instance using R?

3. What is the relationship between the privacy and the utility of synthetically generated data?

*"Differential privacy is a mathematical property of synthetic data generation methods, and rigorously showing that a synthesizer satisfies differential privacy needs thorough mathematical reasoning. Not all differentially private methods are introduced in a rigorous setting, and neither are they built on the same mathematical terminology. This leads to the following question. How are differentially private methods defined and explained using a consequent and precise mathematical language? There are several methods for synthesization of differentially private data such as the smoothed histogram mechanism and PrivBayes, but they are quite theoretical. How can these mechanisms, if possible, be implemented in practice using a programming language such as R? Furthermore, there is an inherent trade-off between the privacy and the utility of synthetically generated data. What is the relationship between the privacy and the utility of synthetic data?"*
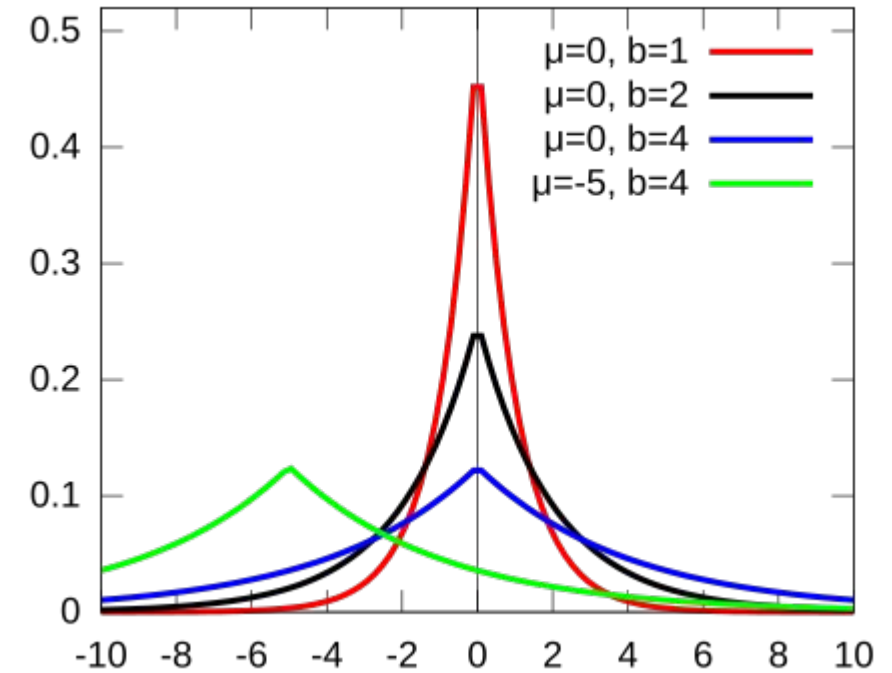
# Differential Privacy

- The relative change of the conditional distribution with respect to a unit change in the original dataset is bounded by *exp(a)*.

- This guarantees a certain level of privacy even in the presence of outliers.

- For *a* close to zero the conditional distributions turn indistinguishable.

$$\sup_{\substack{B \in \mathcal{B}(\mathbb{R}^{k \times d}) \\ (x, x') \in N(X)}} \frac{Q(B|X = x)}{Q(B|X = x')} \leq \exp(\alpha)$$

$$\frac{Q(B|X = x)}{Q(B|X = x')} \in \left[ \exp(-\alpha), \exp(\alpha) \right]$$

# The Laplace mechanism and the exponential mechanism

- The Laplace mechanism:
  - Add Laplace distributed noise to each observation.
  - The scale parameter depends on the sensitivity and the wanted level of differential privacy.
  - Method can only be applied to continuous variables

- The Exponential mechanism:
  - The function $q$ scores a potential output to an input, and outputs are then sampled such that high scores are exponentially more likely to be picked than lower scores.
  - Can be applied to any type of data, but can be difficult to implement.
  - Postulated in the literature to be a foundational differentially private mechanism, meaning that **all differentially private mechanisms are variants of the exponential mechanism.**



$$g_x(r) \propto \exp\left(\alpha q(x,r)\right),$$

$$Q(B|X = x) := \int_B g_x(r)d\mu(r)$$

AALBORG UNIVERSITY

# PrivBayes

- Utilizes Bayesian networks.

- We can implement this using the exponential mechanism for generating Bayesian networks differentially private.

- The Laplace mechanism can be used for injecting noise into the conditional probabilities.

**Algorithm 4.0.1. (PrivBayes)**

Let $X : \Omega \to \mathbb{R}^{n \times d}$ with $n, d \in \mathbb{N}$ be a database, let $V = (X_1, X_2, \ldots, X_d)$ be a random vector that is equal in distribution to each row of $X$. The `PrivBayes` mechanism is given by the following algorithm

1. Fit a $k$-degree Bayesian network

$$\mathcal{N} = \left\{ (1, \Pi_1), (2, \Pi_2), \ldots, (d, \Pi_d) \right\},$$

that resembles the conditional independence relations in $V$ for some chosen $k \in \mathbb{N}$, using an $\alpha_1$-differentially private method.

2. Estimate conditional probabilities $\mathbb{P}(X_i | X_j, j \in \Pi_i)$ for all $i = 1, 2, \ldots, d$, and inject noise using an $\alpha_2$-differentially private method.

3. Using the noisy conditional distributions, assemble a noisy joint distribution, and sample a synthetic dataset.

# Generating some data

$$\begin{bmatrix} 1 & 0.6 & 0.5 & 0.4 & 0.3 \\ 0.6 & 1 & 0.6 & 0.5 & 0.4 \\ 0.5 & 0.6 & 1 & 0.6 & 0.5 \\ 0.4 & 0.5 & 0.6 & 1 & 0.6 \\ 0.3 & 0.4 & 0.5 & 0.6 & 1 \end{bmatrix}$$

$\longrightarrow$

| X1 | X2 | X3 | X4 | X5 |
|----|----|----|----|----|
| 0 | 1 | 1 | 1 | 1 |
| 0 | 1 | 0 | 1 | 1 |
| 0 | 1 | 1 | 1 | 0 |
| 1 | 1 | 0 | 0 | 0 |
| 0 | 1 | 1 | 1 | 1 |
| 0 | 1 | 1 | 1 | 1 |
| 0 | 1 | 0 | 1 | 1 |
| 0 | 0 | 0 | 0 | 0 |
| 1 | 1 | 1 | 1 | 0 |
| 1 | 1 | 0 | 1 | 1 |

$\longrightarrow$

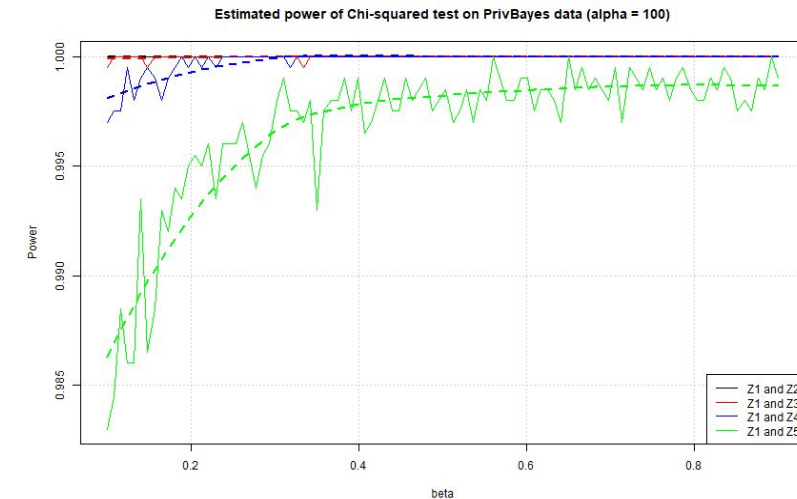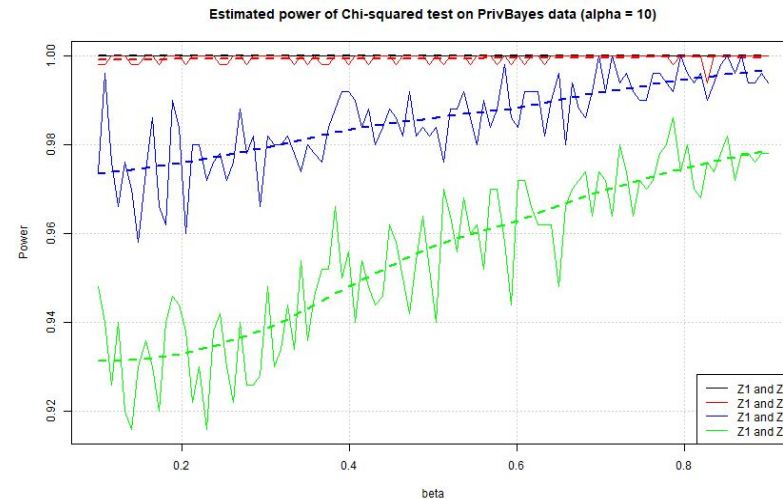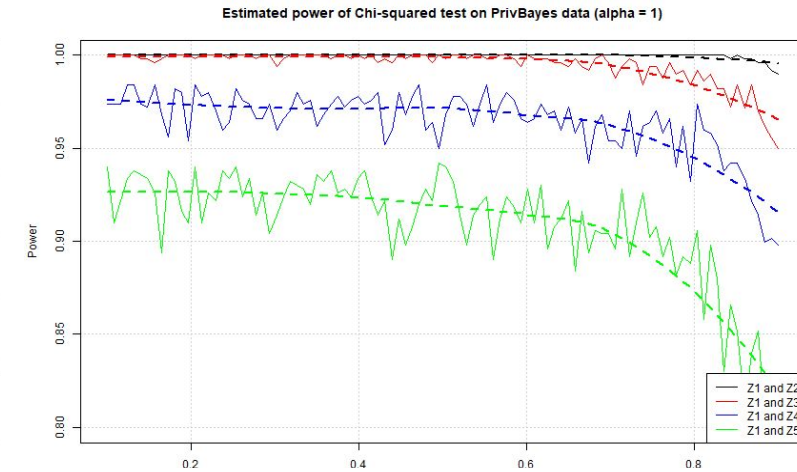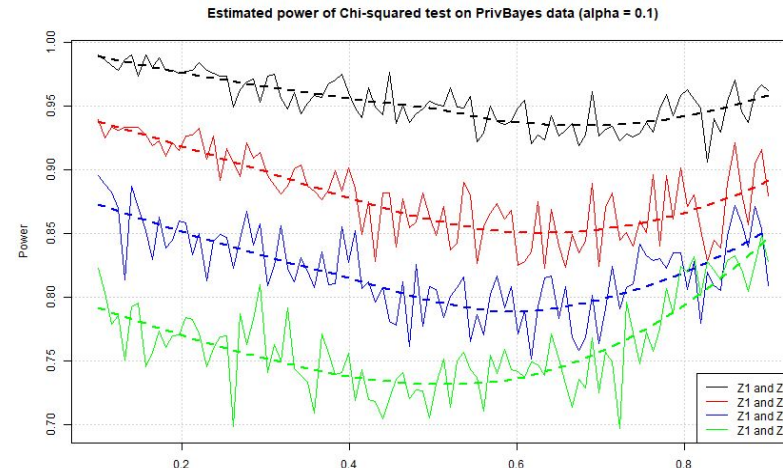| Z1 | Z2 | Z3 | Z4 | Z5 |
|----|----|----|----|----|
| 1 | 1 | 0 | 0 | 0 |
| 0 | 1 | 0 | 1 | 0 |
| 1 | 1 | 0 | 1 | 1 |
| 1 | 1 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 |
| 1 | 1 | 1 | 0 | 1 |
| 0 | 1 | 0 | 1 | 1 |
| 1 | 1 | 0 | 0 | 1 |
| 0 | 1 | 1 | 1 | 0 |
| 1 | 0 | 1 | 1 | 0 |

# The mixture parameter
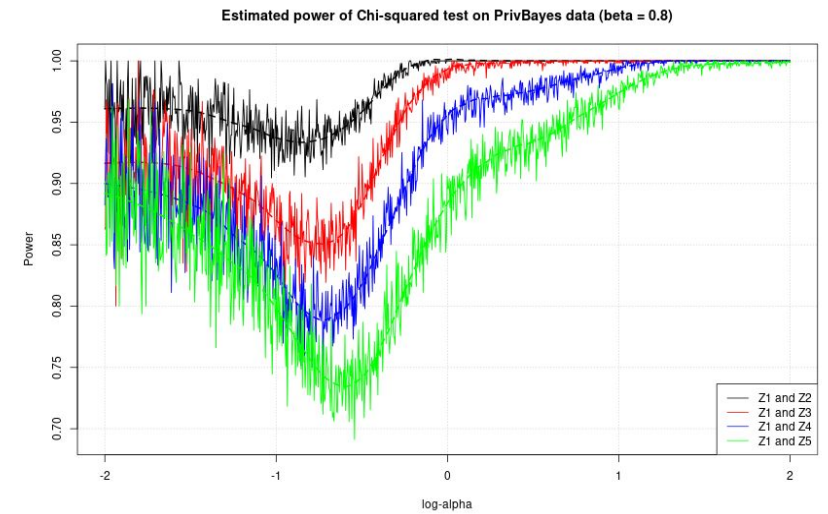
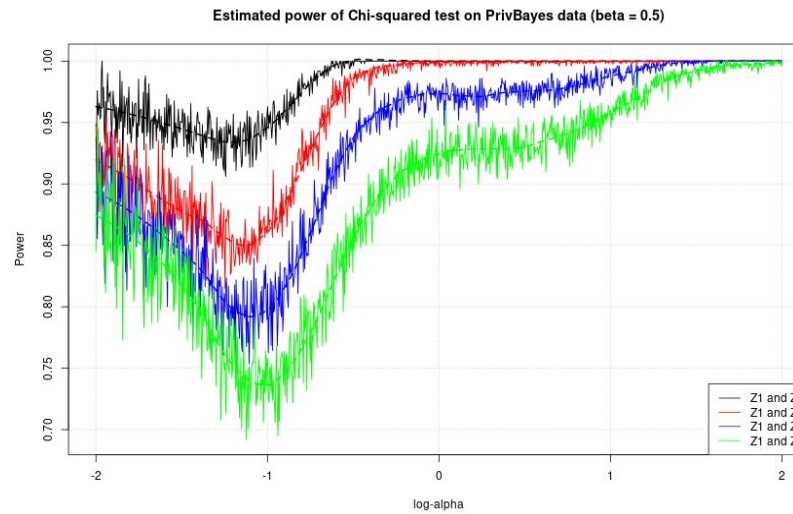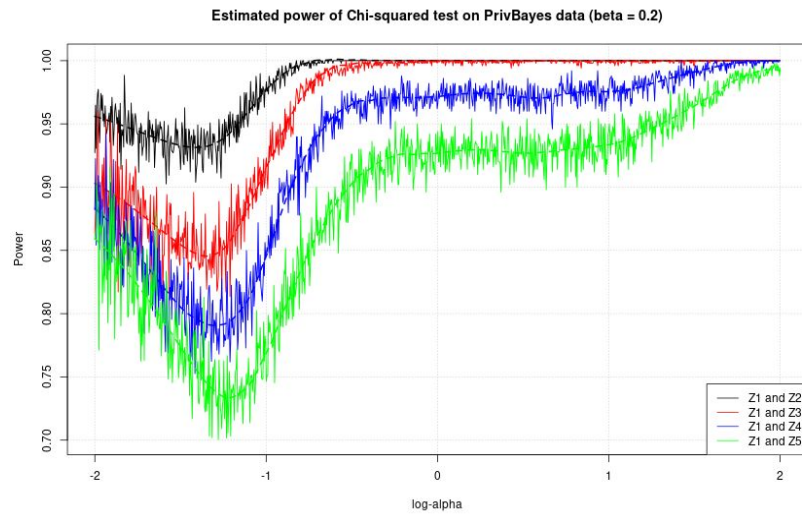$$\alpha := (\alpha_1 + \alpha_2)$$

$$\beta := \frac{\alpha_1}{\alpha} \qquad \beta \in [0, 1]$$

$$\alpha_1 = \beta\alpha$$

$$\alpha_2 = (1 - \beta)\alpha$$

# Overall privacy and power of Pearson's test



Estimated power of Chi-squared test on PrivBayes data (beta = 0.2)

Estimated power of Chi-squared test on PrivBayes data (beta = 0.5)

Estimated power of Chi-squared test on PrivBayes data (beta = 0.8)

# An example

- Dichotomised a diabetes dataset.

- Outcome is 1 for patients who received diabetes diagnosis.

- Will perform a logistic regression for assessing significant variables for getting diabetes diagnosis.

- *Lets go to R!*



| | Outcome | ageUnder45 | overweight | everPregnant |
|---|---|---|---|---|
| [1,] | 1 | 0 | 1 | 1 |
| [2,] | 0 | 1 | 1 | 1 |
| [3,] | 1 | 1 | 0 | 1 |
| [4,] | 0 | 1 | 1 | 1 |
| [5,] | 1 | 1 | 1 | 0 |
| [6,] | 0 | 1 | 1 | 1 |
| [7,] | 1 | 1 | 1 | 1 |
| [8,] | 0 | 1 | 1 | 1 |
| [9,] | 1 | 0 | 1 | 1 |
| [10,] | 1 | 0 | 0 | 1 |
| [11,] | 0 | 1 | 1 | 1 |
| [12,] | 1 | 1 | 1 | 1 |
| [13,] | 0 | 0 | 1 | 1 |
| [14,] | 1 | 0 | 1 | 1 |
| [15,] | 1 | 0 | 1 | 1 |
| [16,] | 1 | 1 | 1 | 1 |
| [17,] | 1 | 1 | 1 | 0 |
| [18,] | 1 | 1 | 1 | 1 |
| [19,] | 0 | 1 | 1 | 1 |
| [20,] | 1 | 1 | 1 | 1 |
| [21,] | 0 | 1 | 1 | 1 |

# Questions