

General Purpose f -DP Estimation and Auditing in a Black-Box Setting

Önder Askin, Holger Dette, Martin Dunsche,
Tim Kutta, Yun Lu, Yu Wei and Vassilis Zikas

RUHR
UNIVERSITÄT
BOCHUM

RUB



University
of Victoria



Georgia Tech College of Computing
School of Cybersecurity
and Privacy

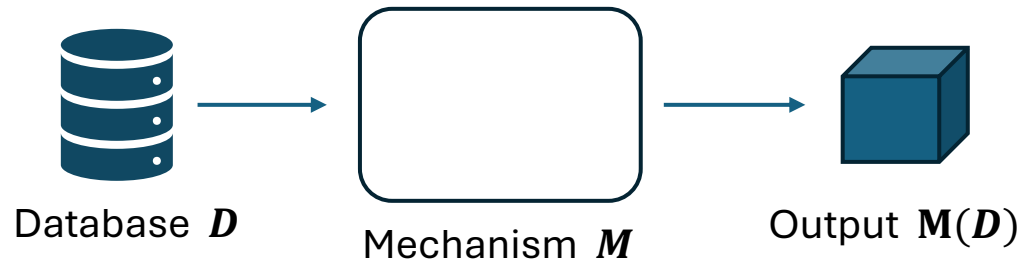
Overview

- f -Differential Privacy
- Estimation
- Auditing

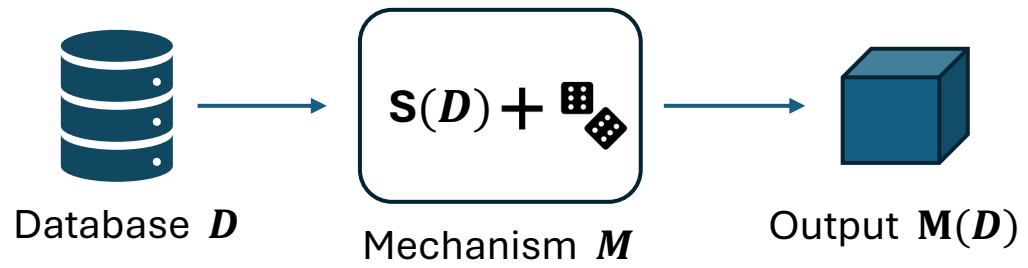
Overview

- f -Differential Privacy
- Estimation
- Auditing

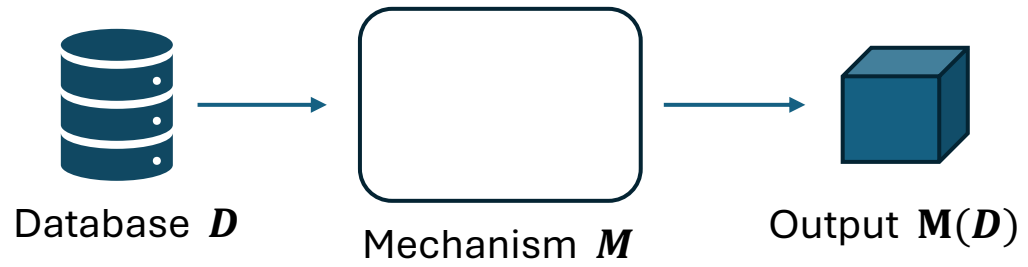
f -Differential Privacy



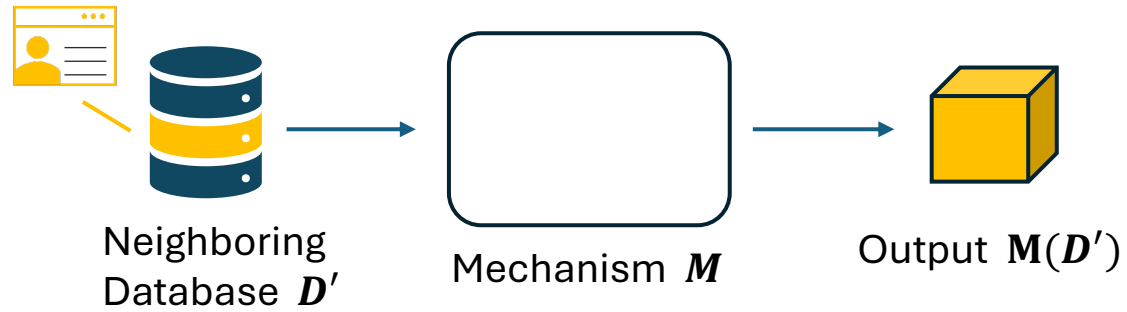
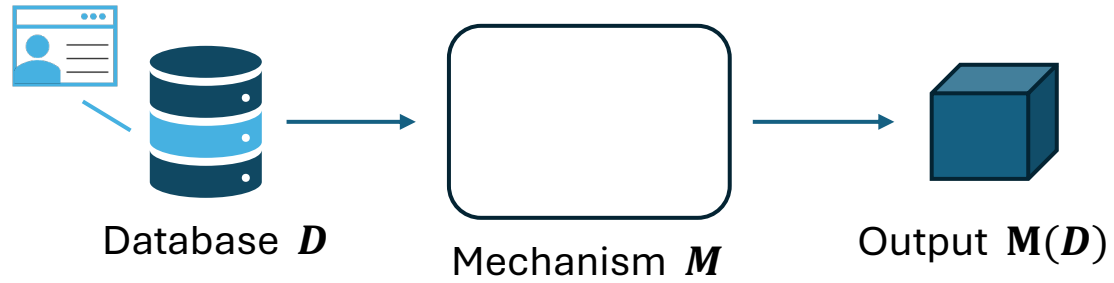
f -Differential Privacy



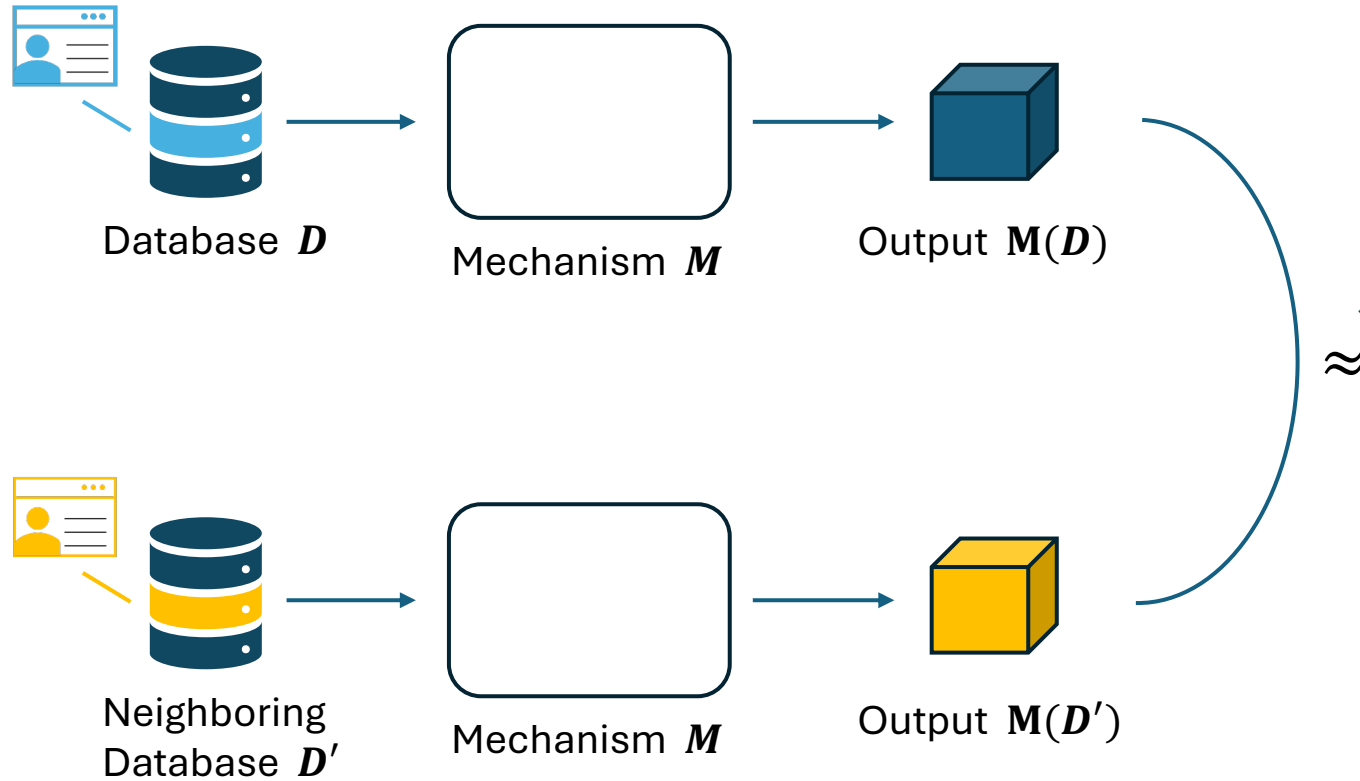
f -Differential Privacy



f -Differential Privacy



f -Differential Privacy



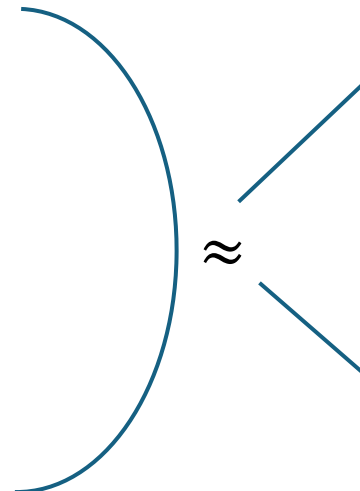
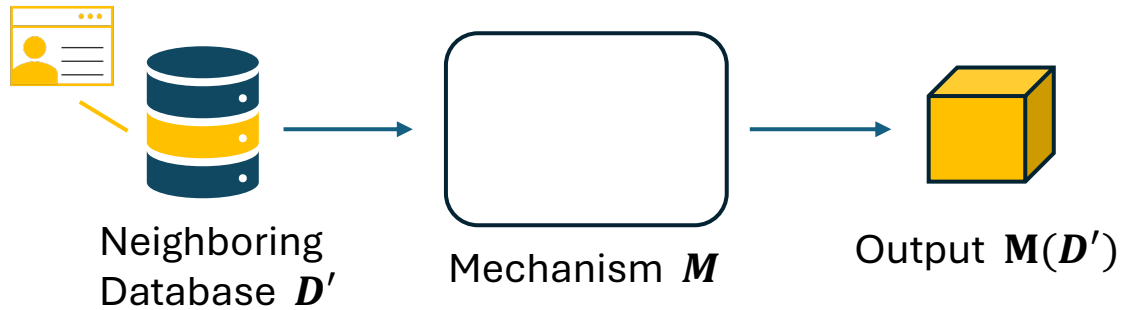
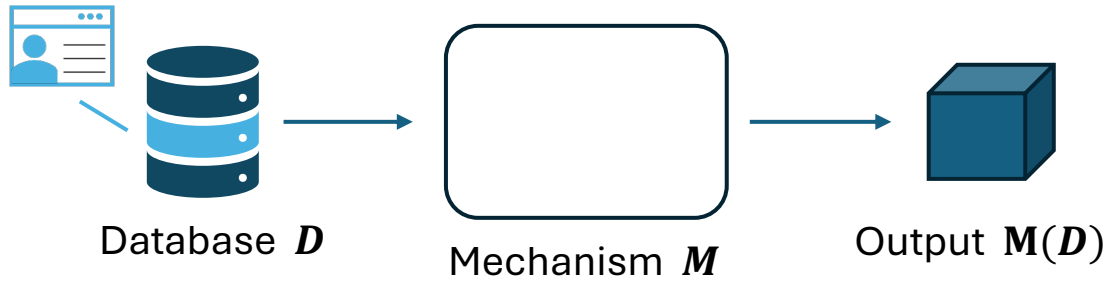
Standard Differential Privacy (DP)

Mechanism M is (ϵ, δ) -DP, if

$$\Pr(M(D) \in S) \leq e^\epsilon \Pr(M(D') \in S) + \delta$$

for all neighboring databases D, D' and sets S

f -Differential Privacy

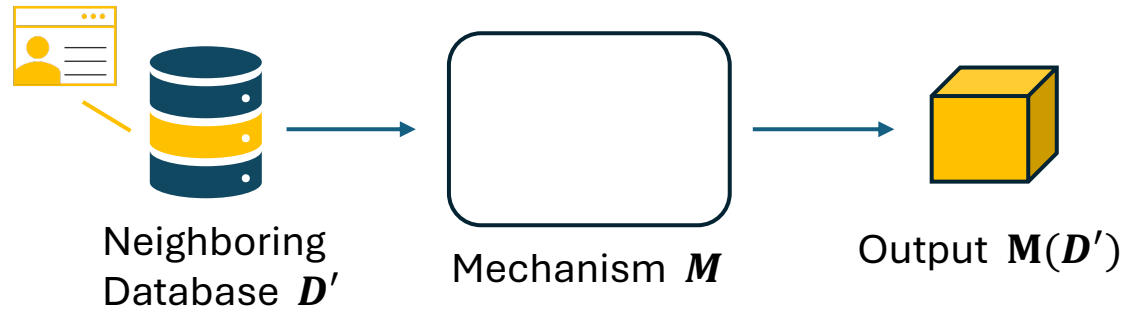
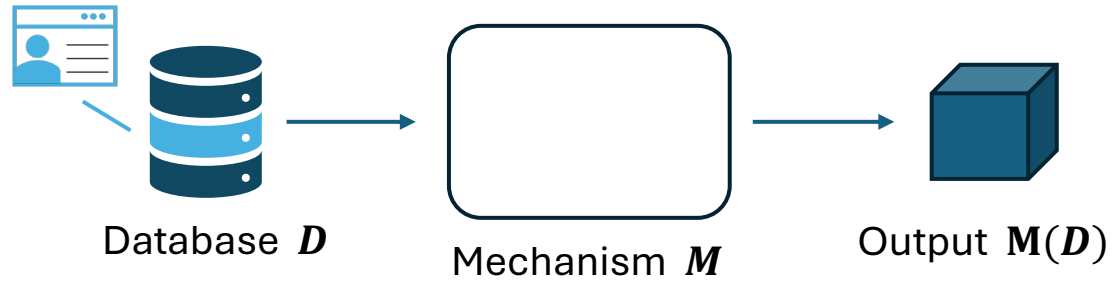


Standard Differential Privacy (DP)
Mechanism M is (ϵ, δ) -DP, if
$$\Pr(M(D) \in S) \leq e^\epsilon \Pr(M(D') \in S) + \delta$$
for all neighboring databases D, D' and sets S

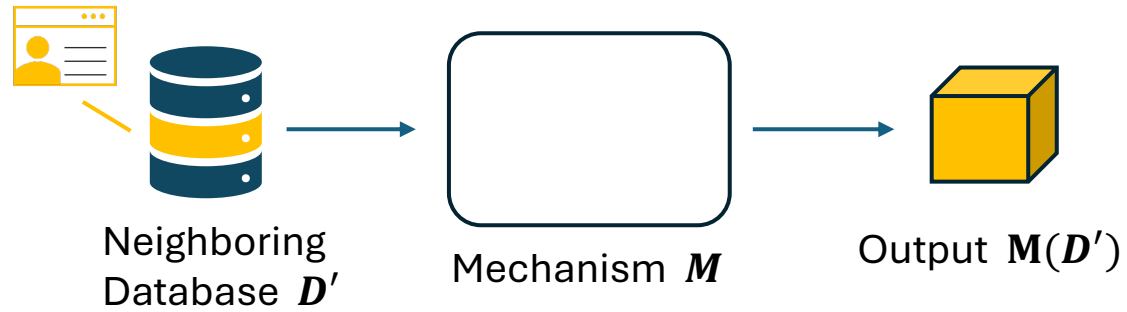
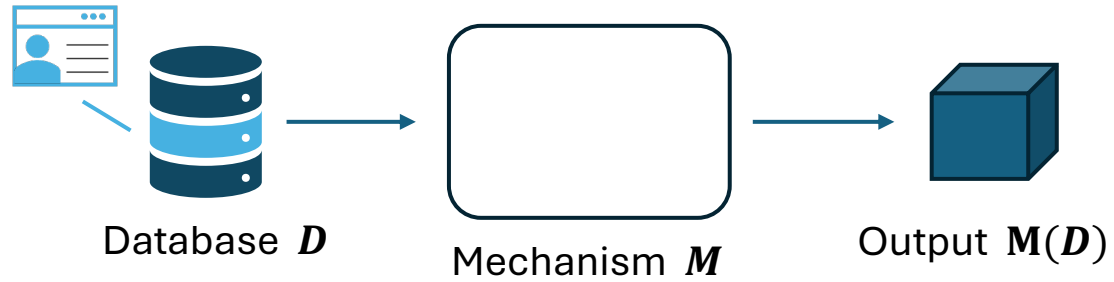
Parameters in Standard DP :
 $\epsilon \in [0, \infty]$

 $\delta \in [0, 1]$

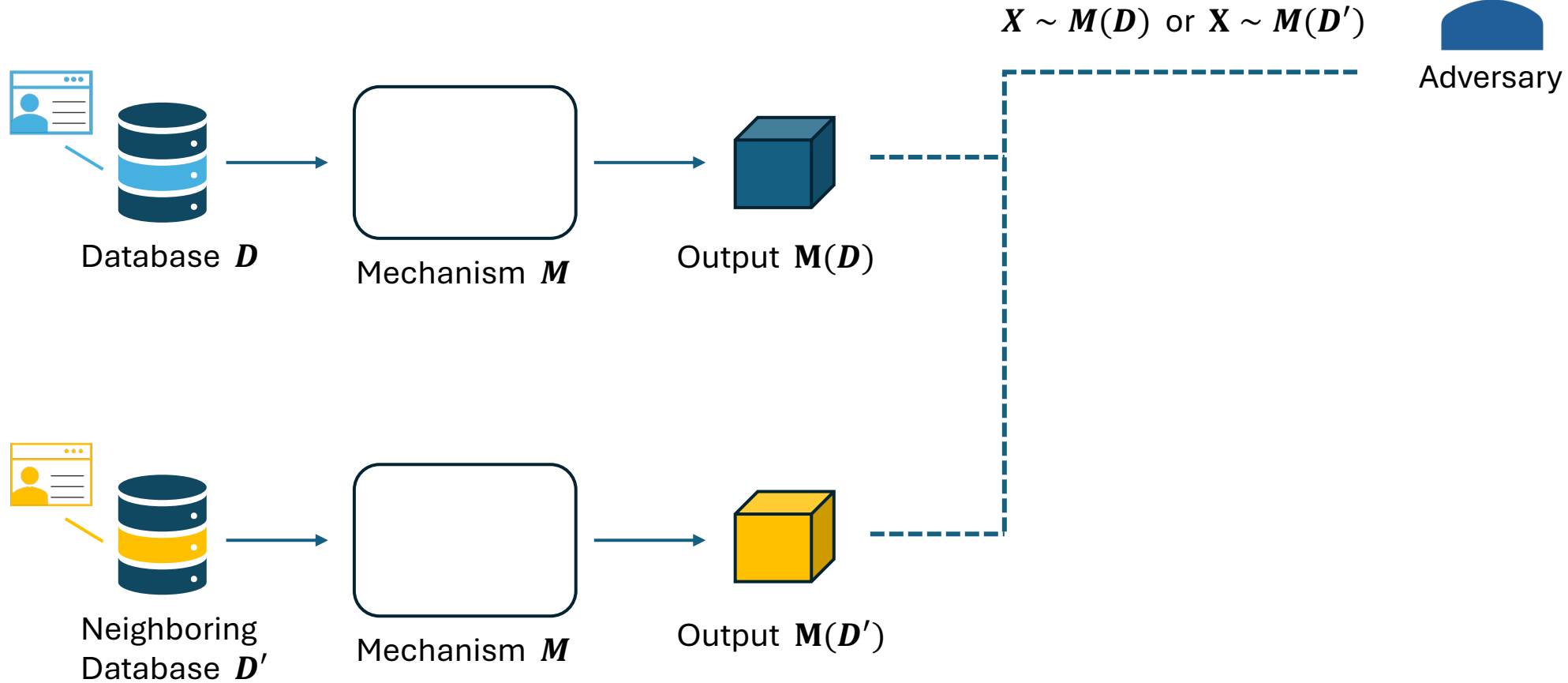
f -Differential Privacy



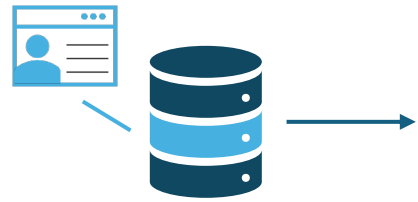
f -Differential Privacy



f -Differential Privacy



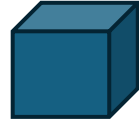
f -Differential Privacy



Database D



Mechanism M



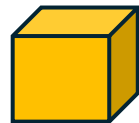
Output $M(D)$



Neighboring
Database D'



Mechanism M



Output $M(D')$

$X \sim M(D)$ or $X \sim M(D')$



Adversary

Problem

$H_0: X \sim M(D)$ or $H_1: X \sim M(D')$

Hypothesis Test

$$g: \mathbb{R}^d \rightarrow \{0, 1\}$$

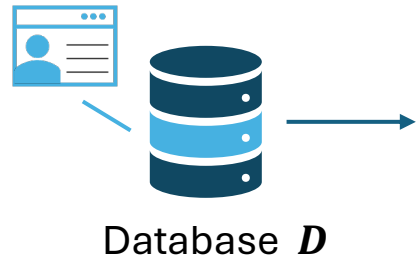
Type-I Error

$$\alpha_g = \Pr_{X \sim M(D)}[g(X) = 1]$$

Type-II Error

$$\beta_g = \Pr_{X \sim M(D')}[g(X) = 0]$$

f -Differential Privacy



$X \sim M(D)$ or $X \sim M(D')$



Adversary

Smallest Type-II Error at Level α

$$\beta = \min \{ \beta_g : g \text{ is a test with } \alpha_g \leq \alpha \}$$

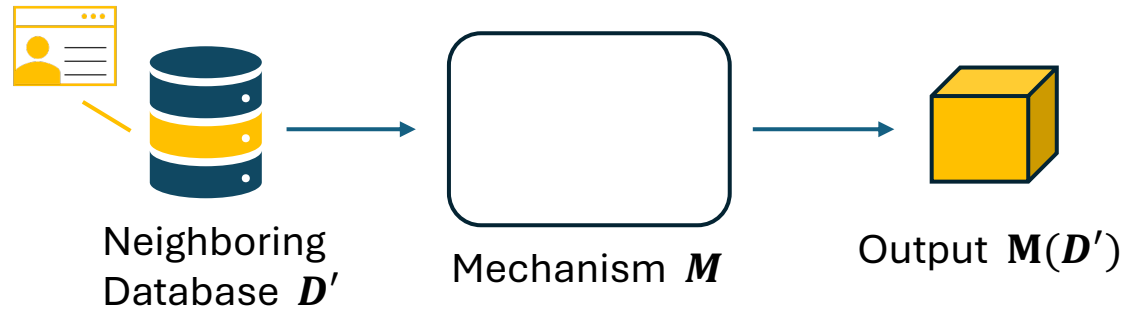
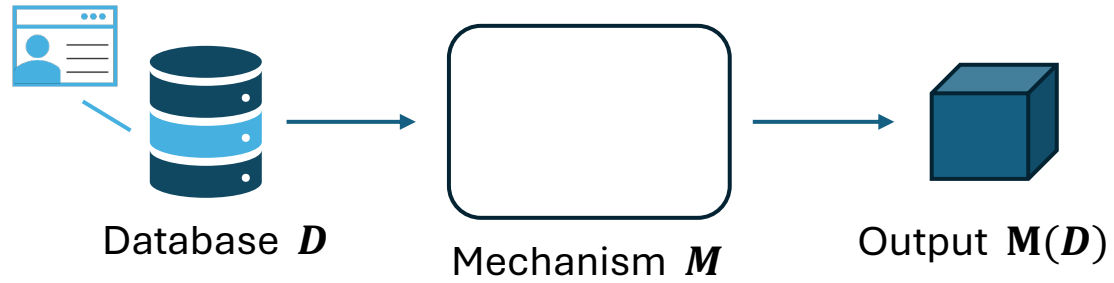
Optimal Test: Likelihood Ratio (LR)

- Probability densities

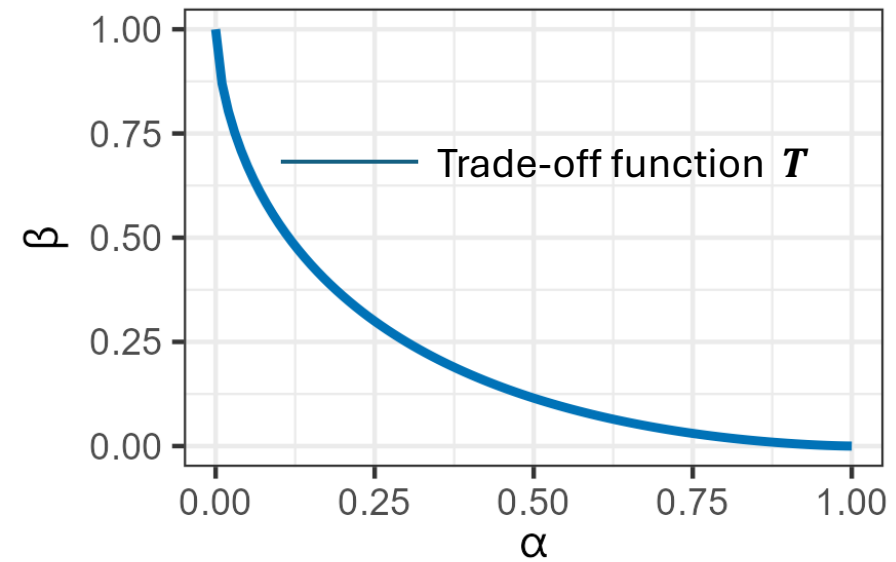
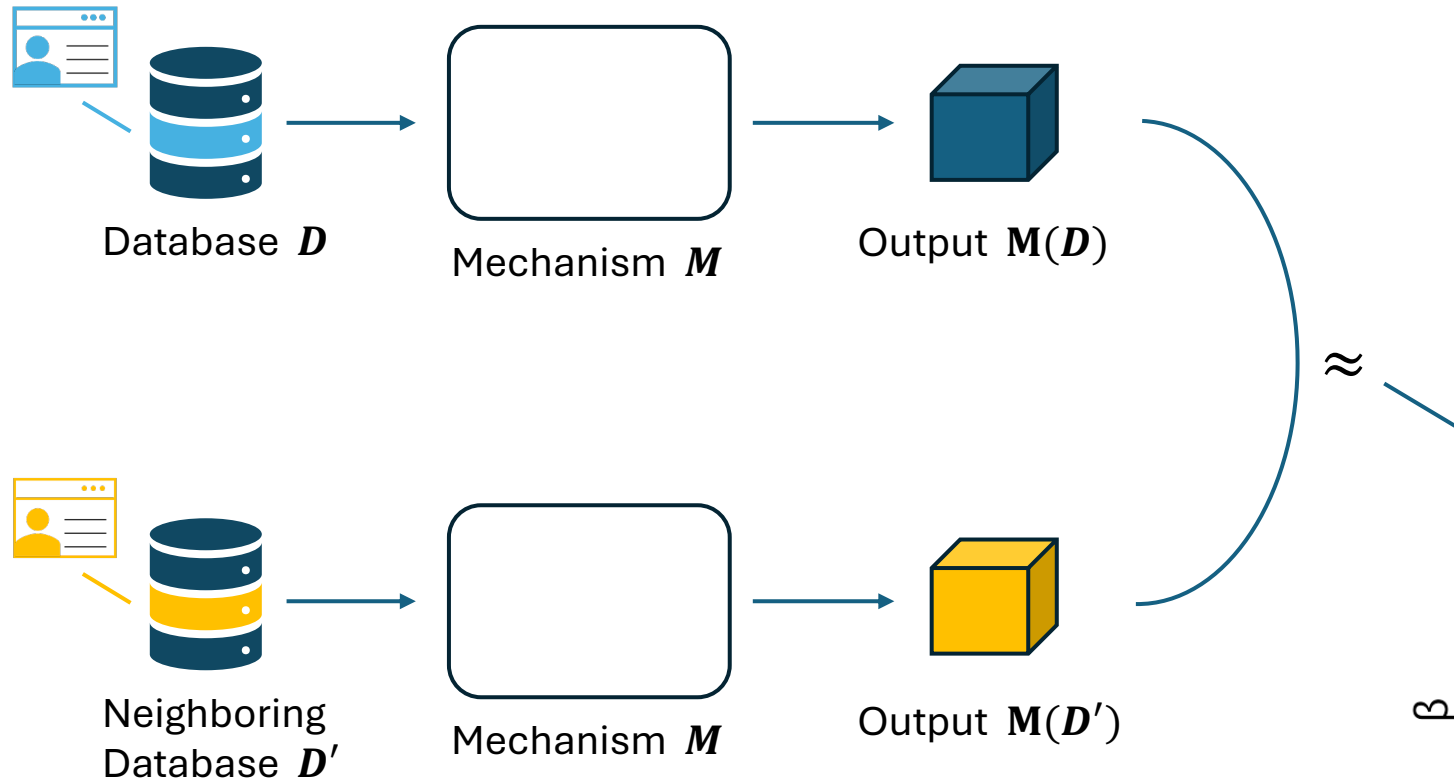
$$\mathbf{p} \sim M(D) \text{ and } \mathbf{q} \sim M(D')$$

- Constants $\eta \geq 0$ and $\lambda \in [0, 1]$
- **Reject** H_0 if $q(X)/p(X) > \eta$
- **Reject** H_0 if $q(X)/p(X) = \eta$ and a coin toss with $\Pr(\text{Heads}) = \lambda$ yields **Heads**

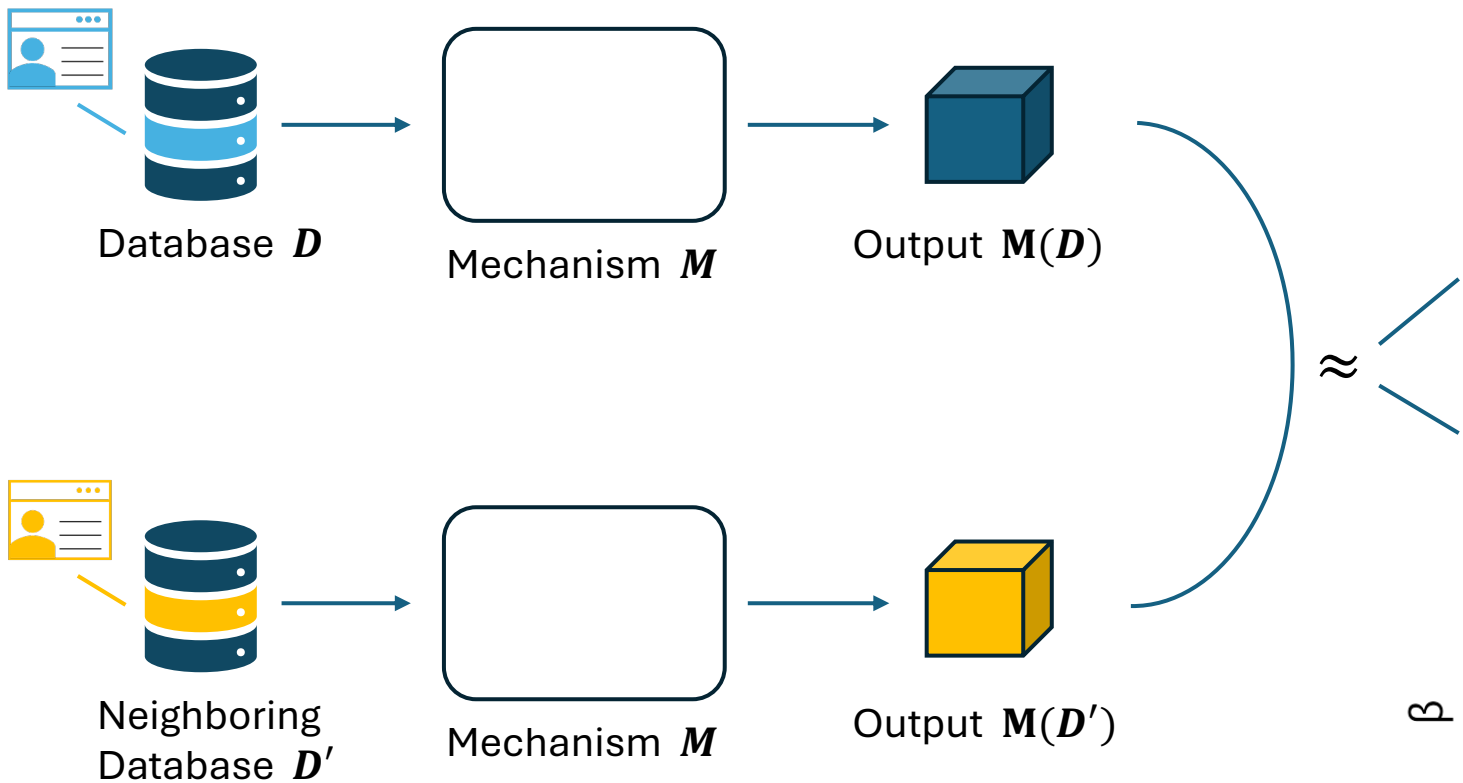
f -Differential Privacy



f -Differential Privacy



f -Differential Privacy

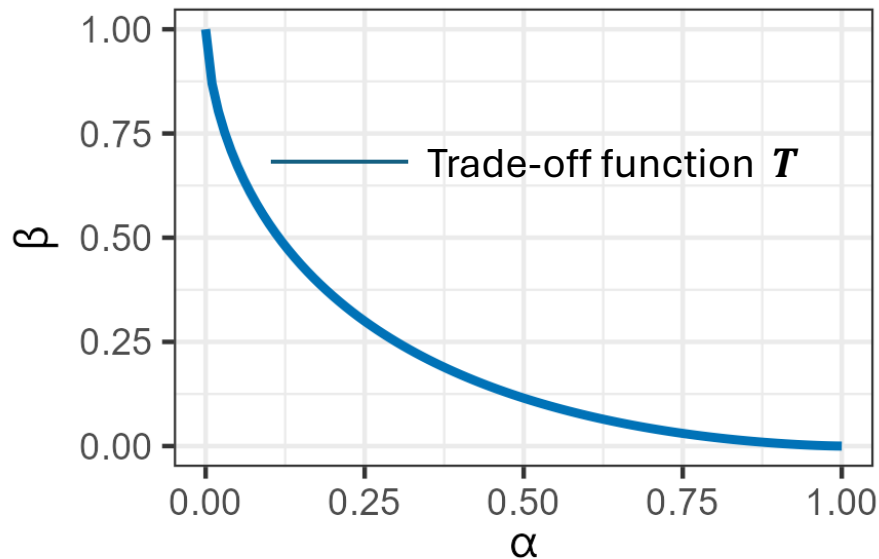


f -Differential Privacy (f -DP)

Mechanism M is f -DP, if for all neighboring databases D, D' we have

$$T(\alpha) \geq f(\alpha) \quad \forall \alpha \in [0,1],$$

where T is implied by $p \sim M(D)$ and $q \sim M(D')$.



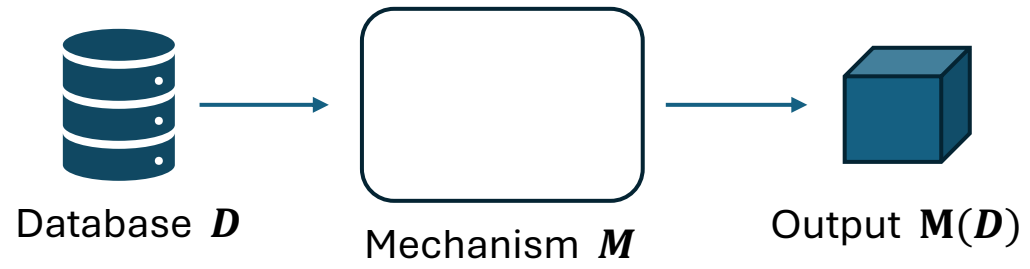
f -Differential Privacy

- comes with easily interpretable privacy guarantees.
- allows for lossless reasoning about composition.
- is suitable for studying privacy amplification.
- is a popular framework for auditing private machine learning.

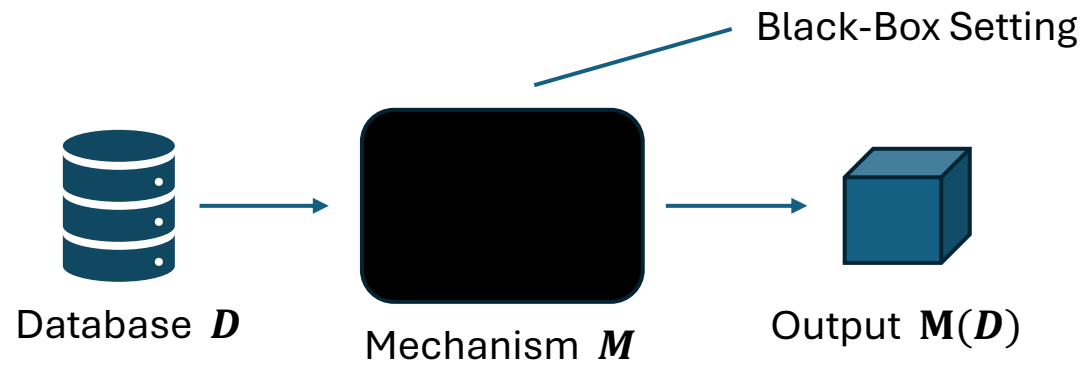
Overview

- f -Differential Privacy
- Estimation
- Auditing

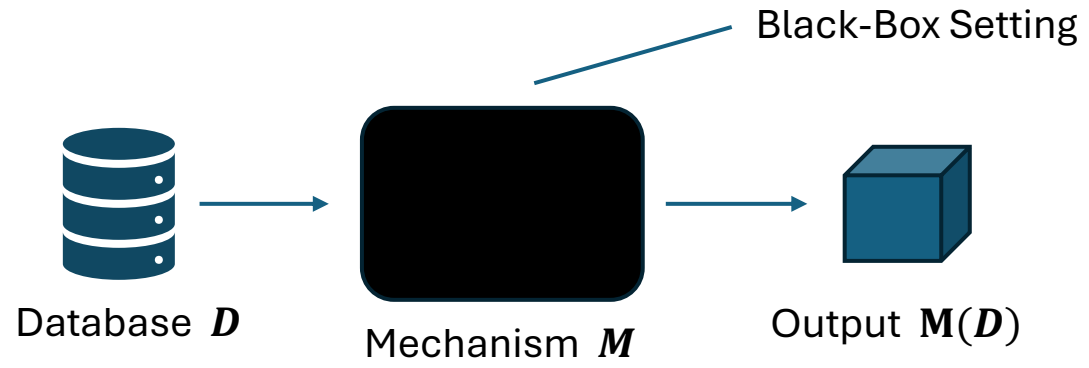
Estimation



Estimation



Estimation



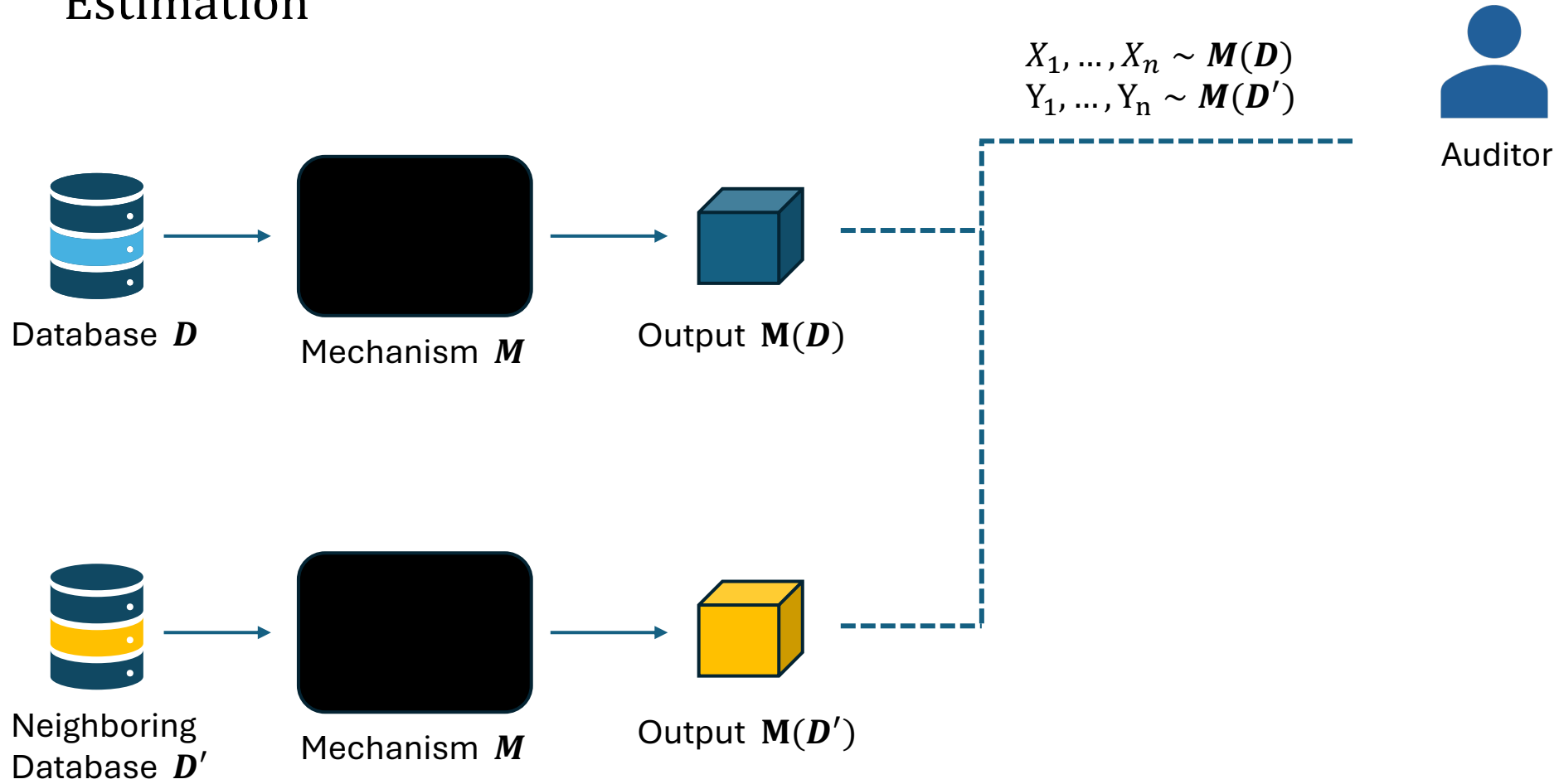
Auditor

Estimation

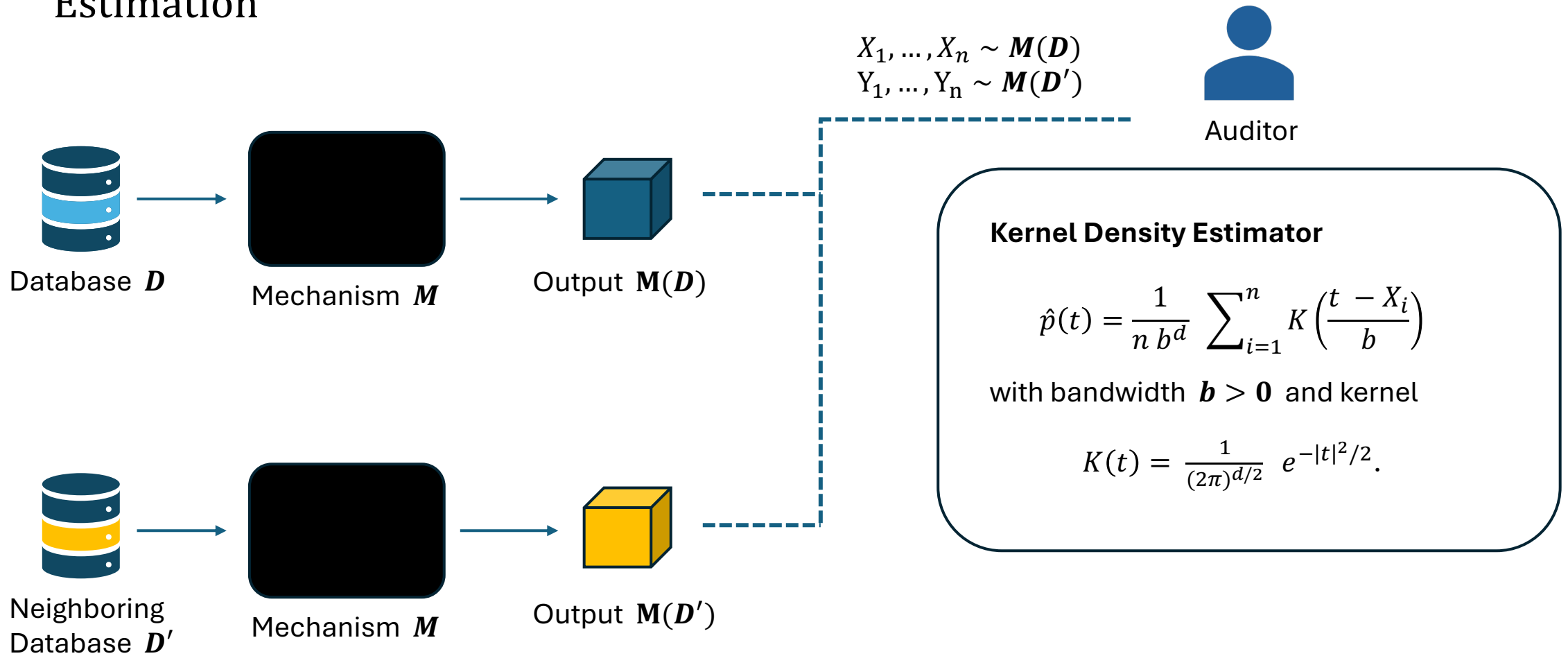


Auditor

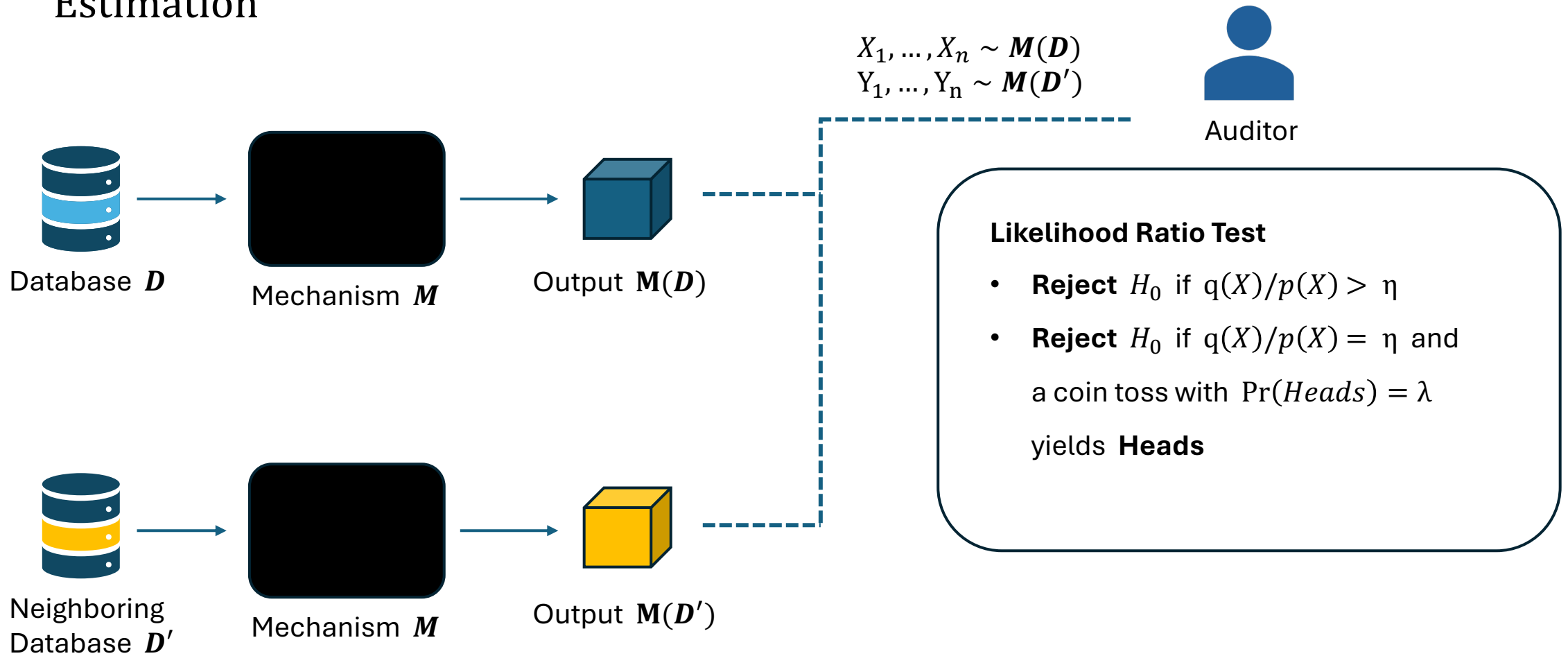
Estimation



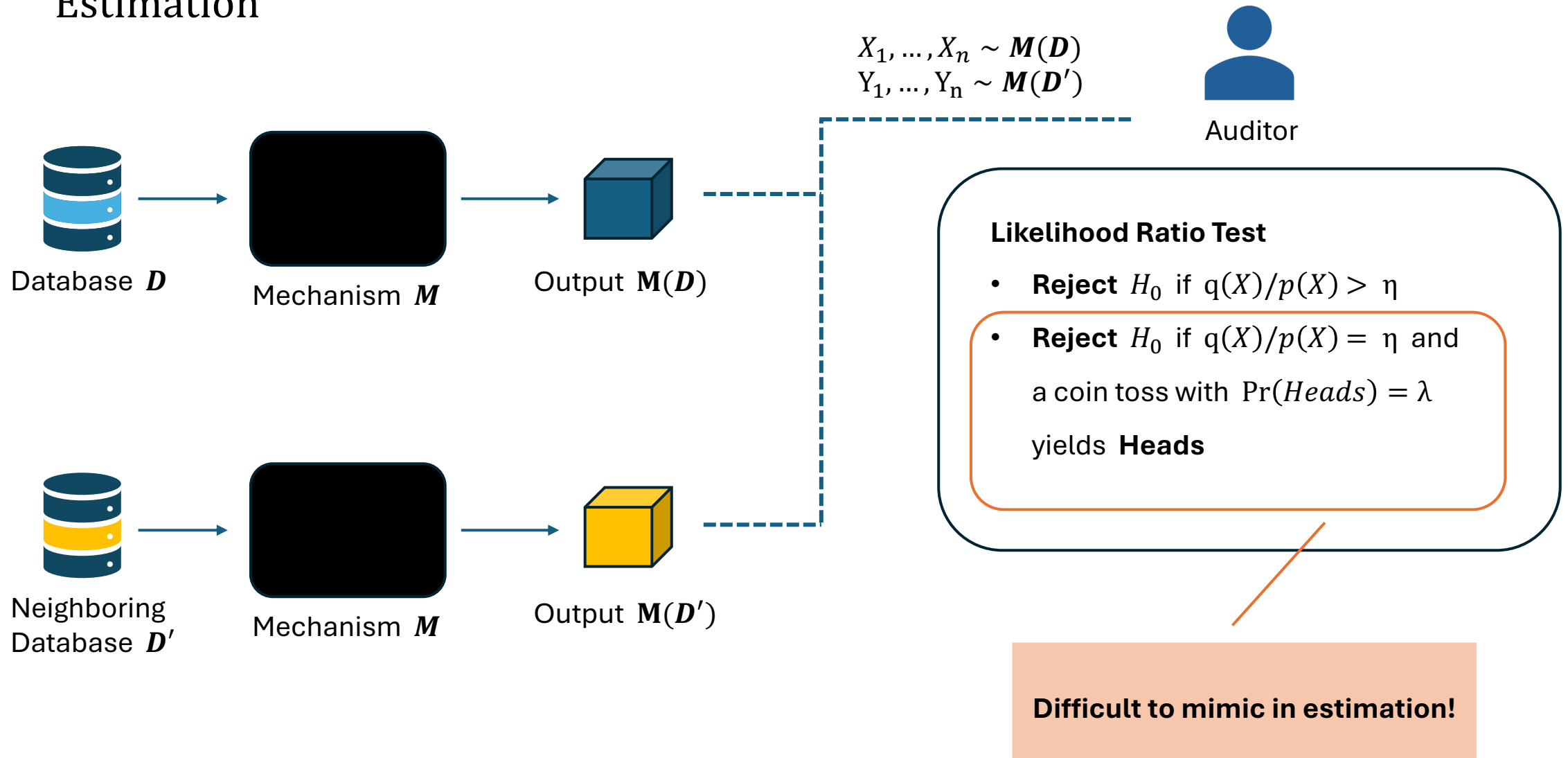
Estimation



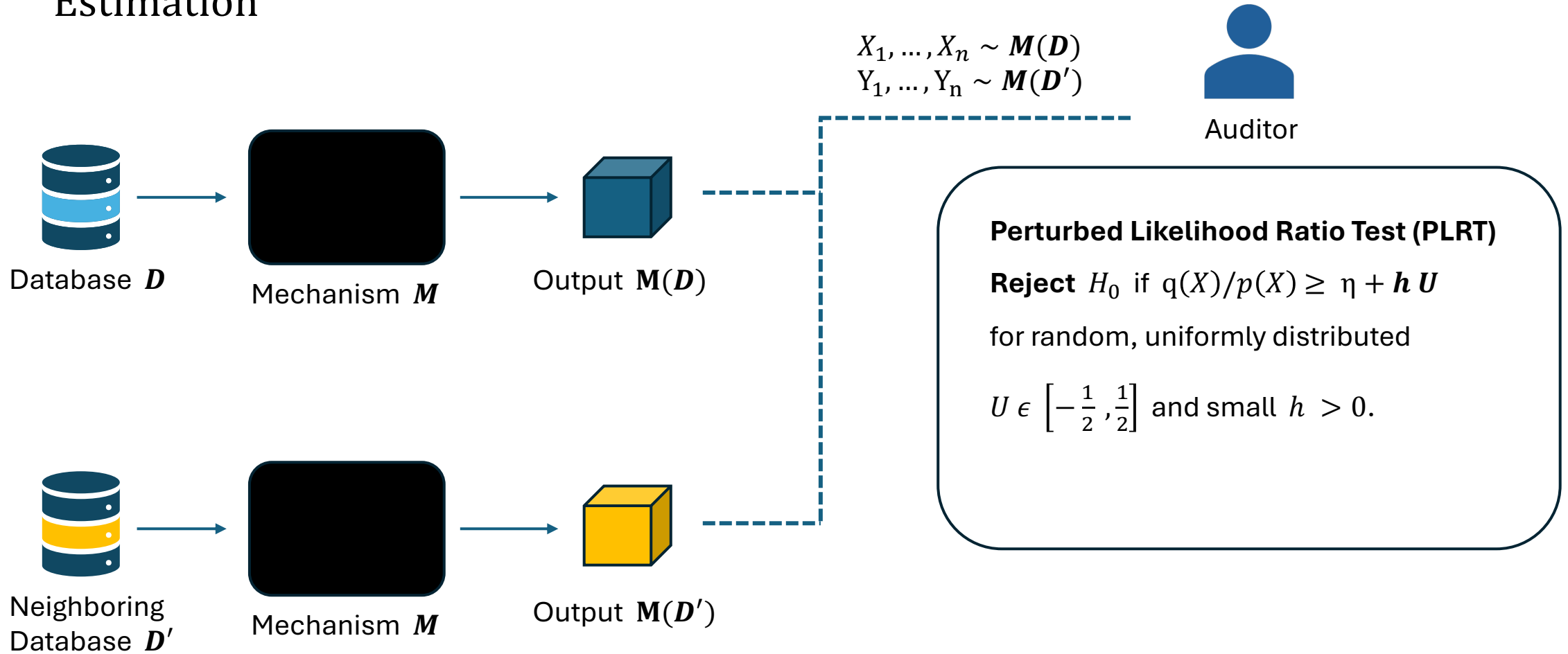
Estimation



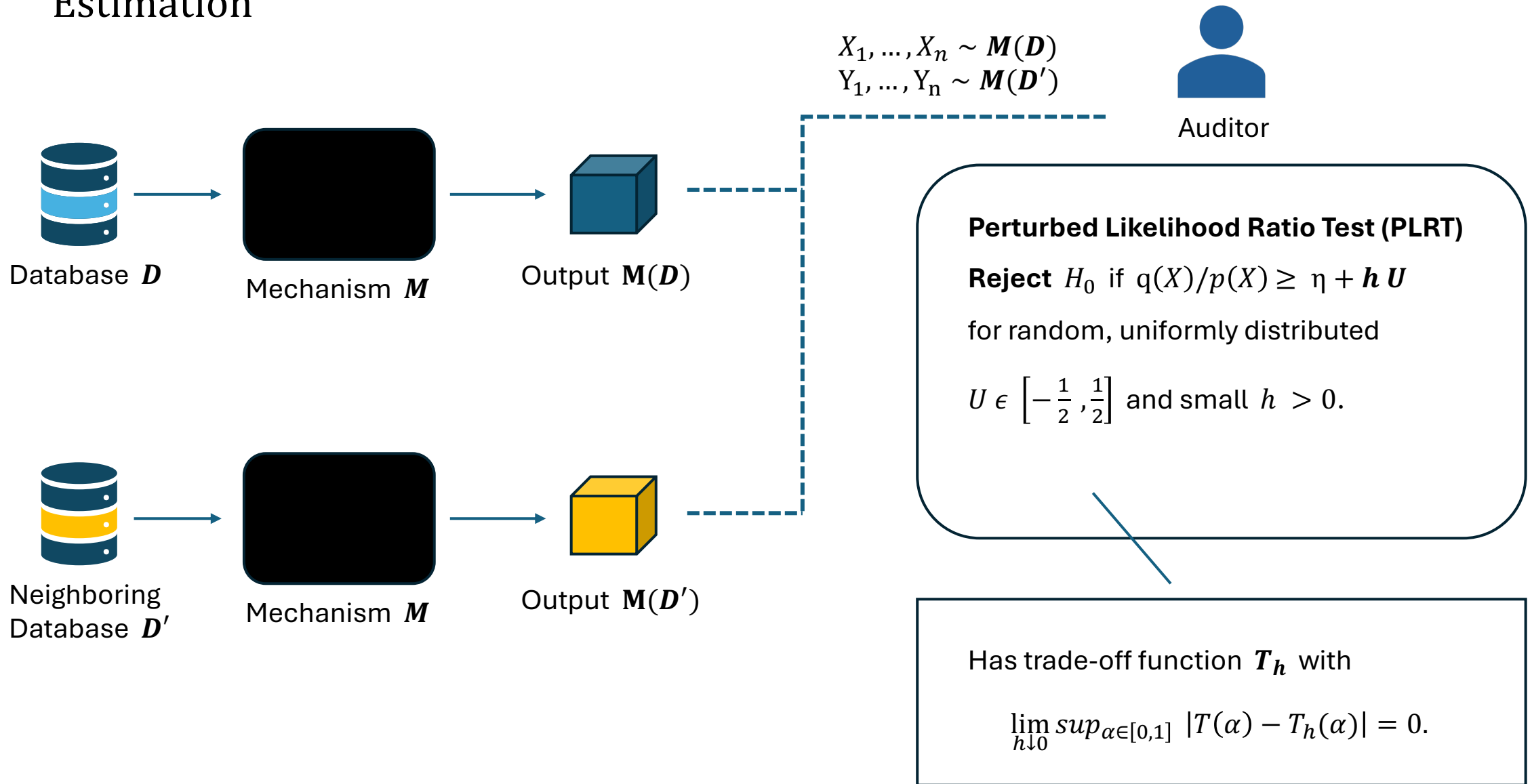
Estimation



Estimation

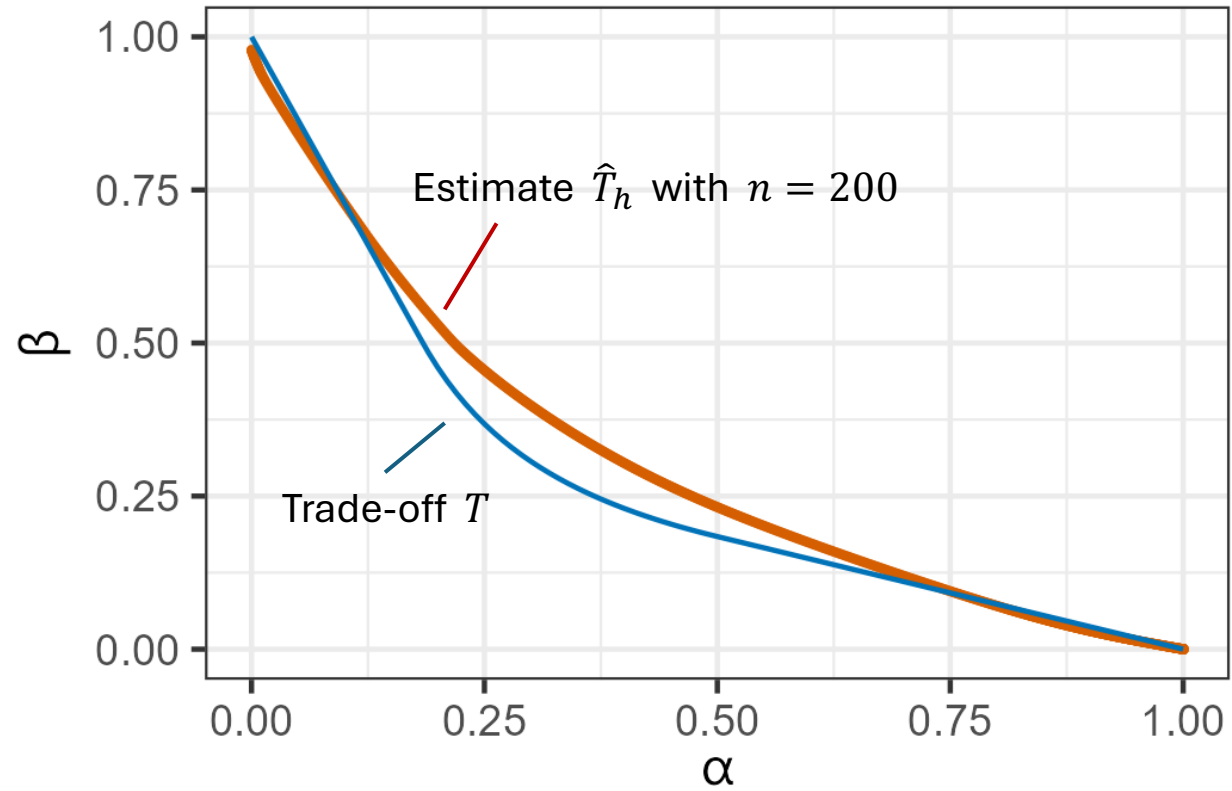


Estimation



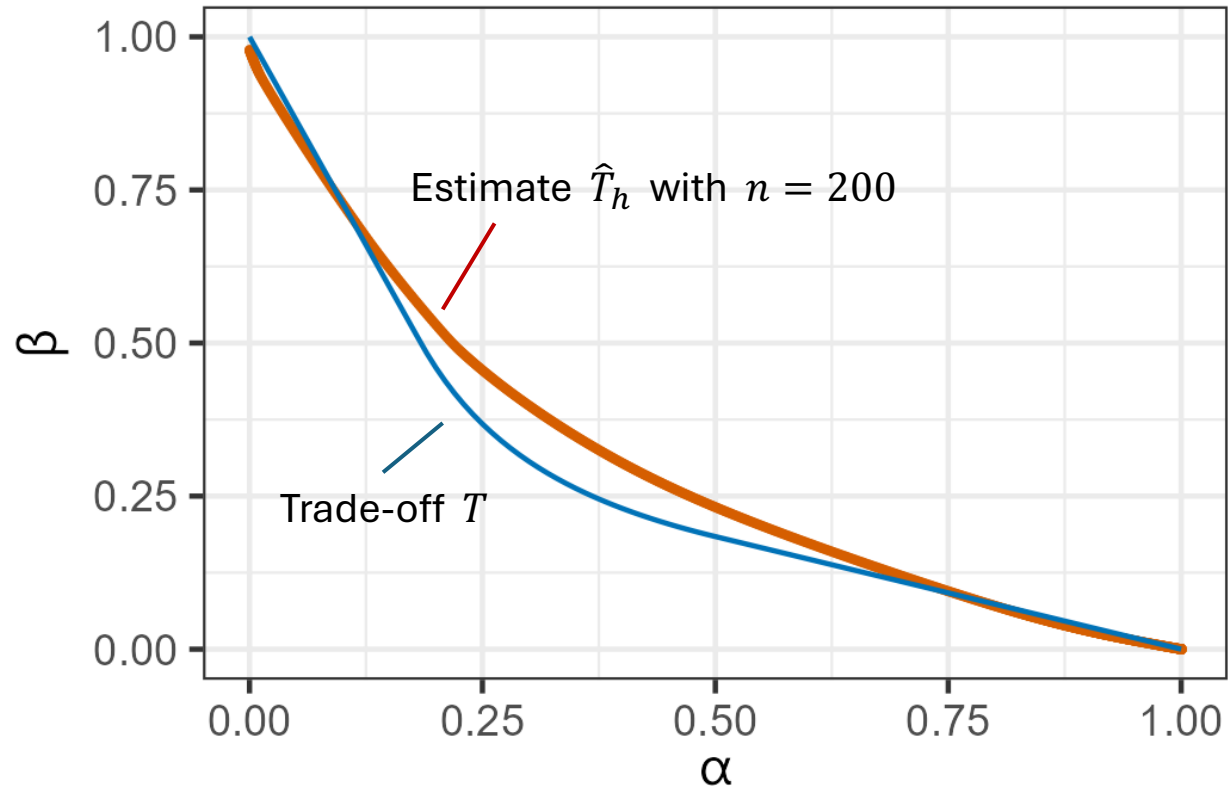
Estimation

Laplace Mechanism



Estimation

Laplace Mechanism



Points of \hat{T}_h

$$\hat{\alpha}_h(\eta) = \int_{x \in [-\frac{h}{2}, \frac{h}{2}]} \frac{1}{h} \int_{\frac{\hat{q}}{\hat{p}} > \eta+x} \hat{p}$$

$$\hat{\beta}_h(\eta) = 1 - \int_{x \in [-\frac{h}{2}, \frac{h}{2}]} \frac{1}{h} \int_{\frac{\hat{q}}{\hat{p}} > \eta+x} \hat{q}$$

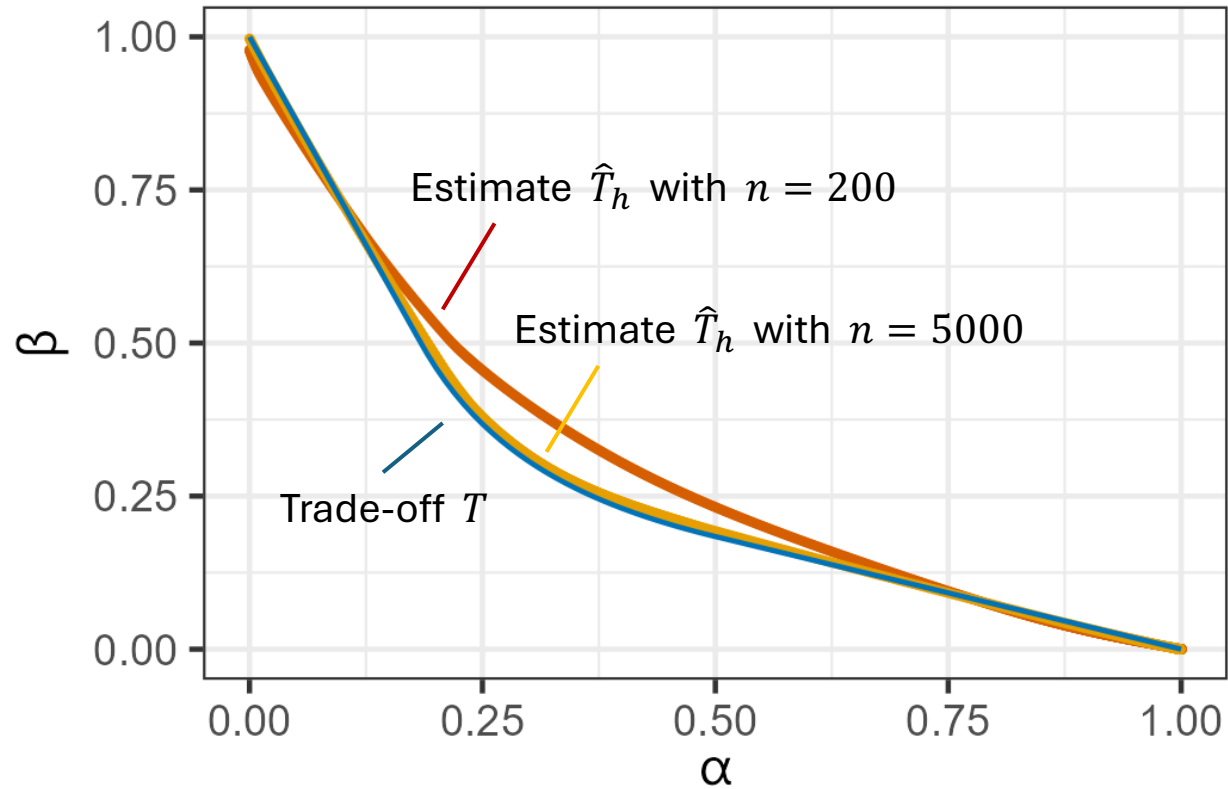
Uniform Convergence

Under certain regularity assumptions, we have for $n \rightarrow \infty$ and $h = h_n \rightarrow 0$ that

$$\sup_{\alpha \in [0,1]} |T(\alpha) - \hat{T}_h(\alpha)| = o_P(1).$$

Estimation

Laplace Mechanism



Points of \hat{T}_h

$$\hat{\alpha}_h(\eta) = \int_{x \in [-\frac{h}{2}, \frac{h}{2}]} \frac{1}{h} \int_{\frac{\hat{q}}{\hat{p}} > \eta+x} \hat{p}$$

$$\hat{\beta}_h(\eta) = 1 - \int_{x \in [-\frac{h}{2}, \frac{h}{2}]} \frac{1}{h} \int_{\frac{\hat{q}}{\hat{p}} > \eta+x} \hat{q}$$

Uniform Convergence

Under certain regularity assumptions, we have for $n \rightarrow \infty$ and $h = h_n \rightarrow 0$ that

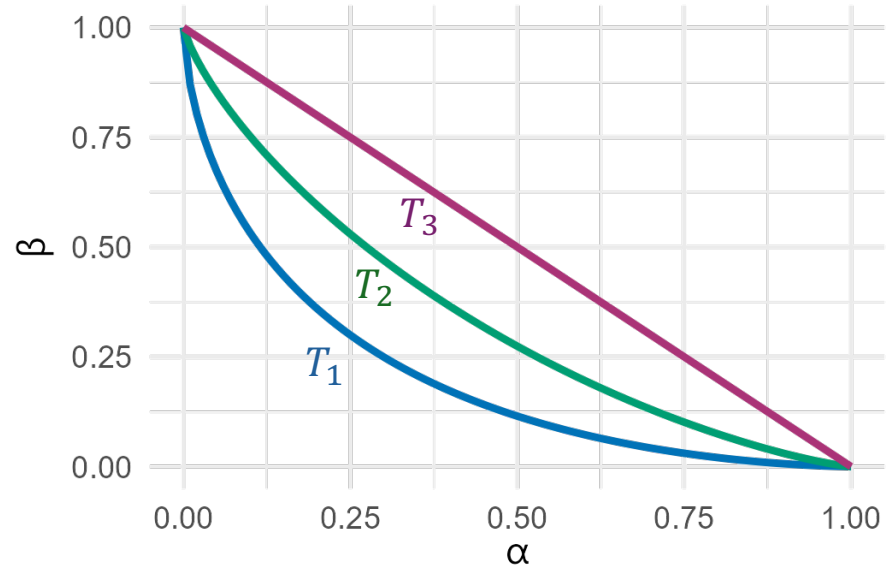
$$\sup_{\alpha \in [0,1]} |T(\alpha) - \hat{T}_h(\alpha)| = o_P(1).$$

Overview

- f -Differential Privacy
- Estimation
- Auditing

Auditing

Privacy Parameters in f -DP

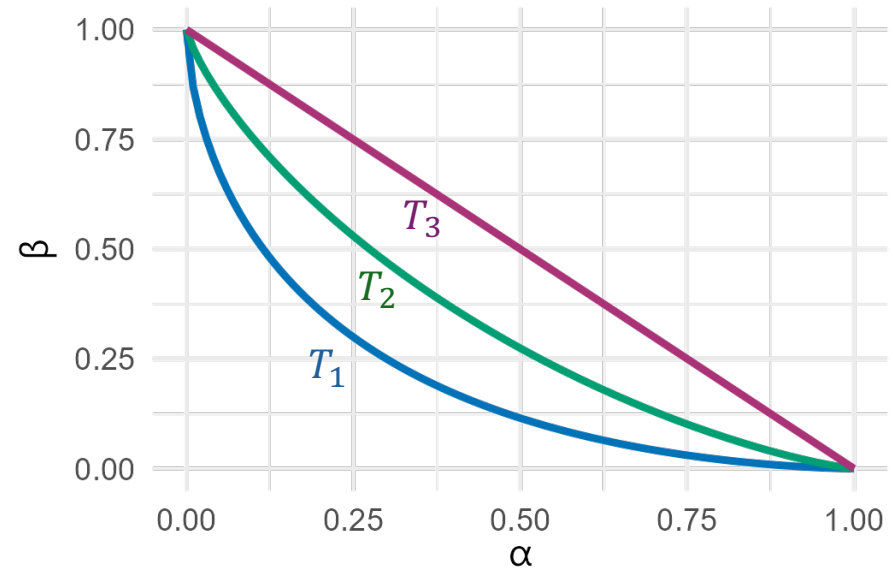


Privacy Parameters in Standard DP

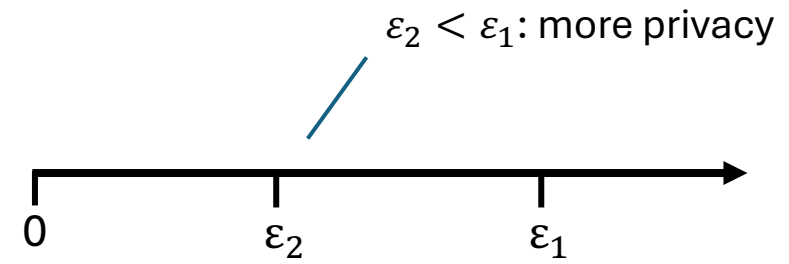


Auditing

Privacy Parameters in f -DP

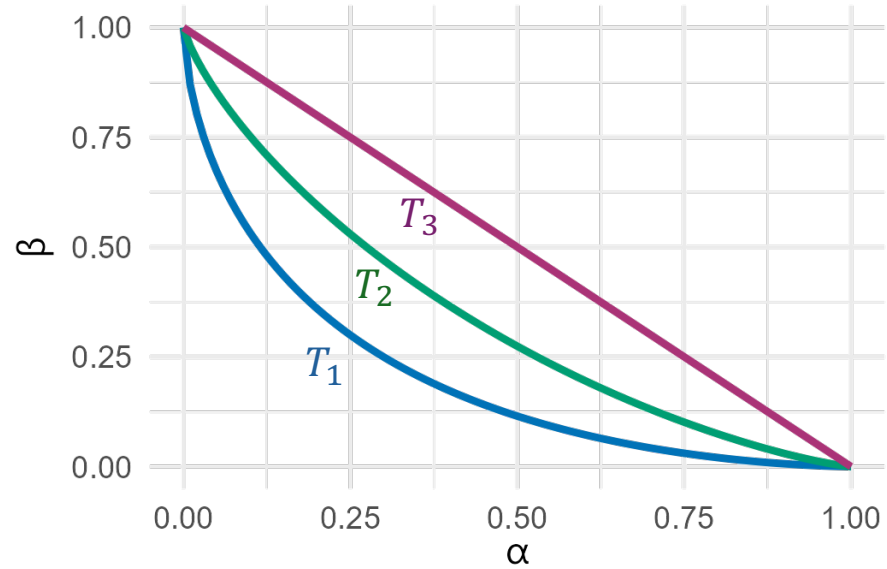


Privacy Parameters in Standard DP

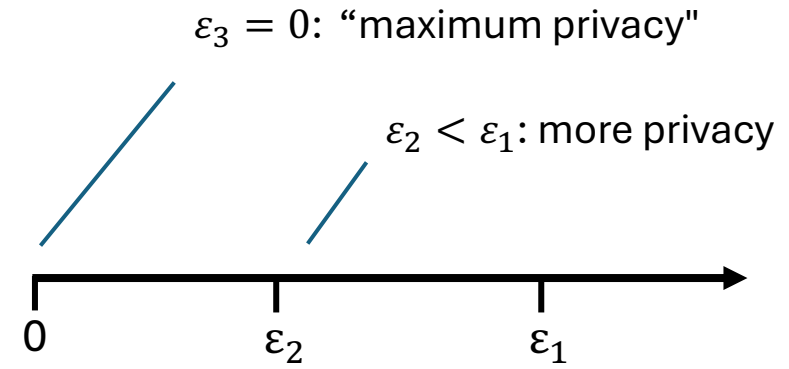


Auditing

Privacy Parameters in f -DP

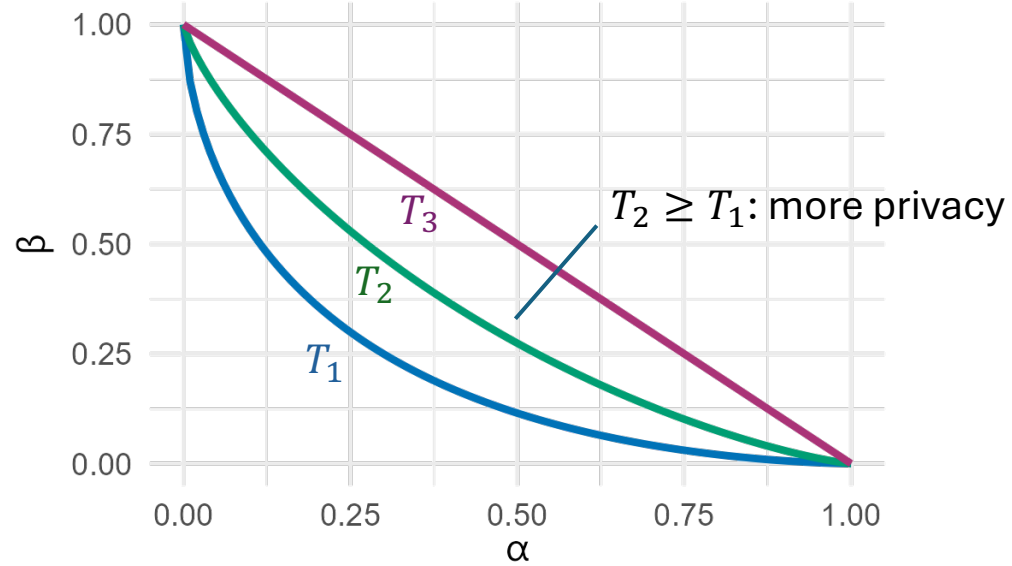


Privacy Parameters in Standard DP

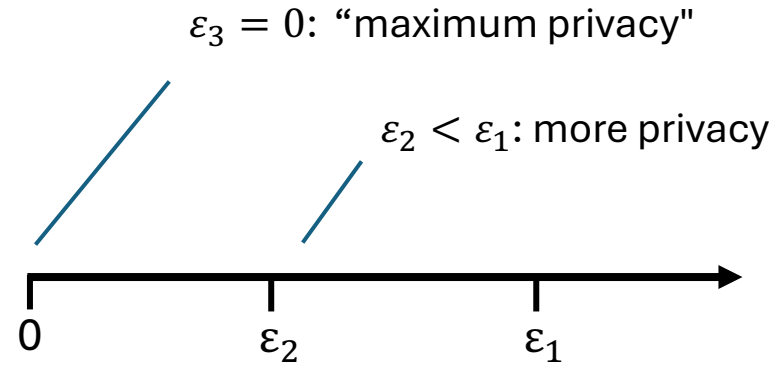


Auditing

Privacy Parameters in f -DP

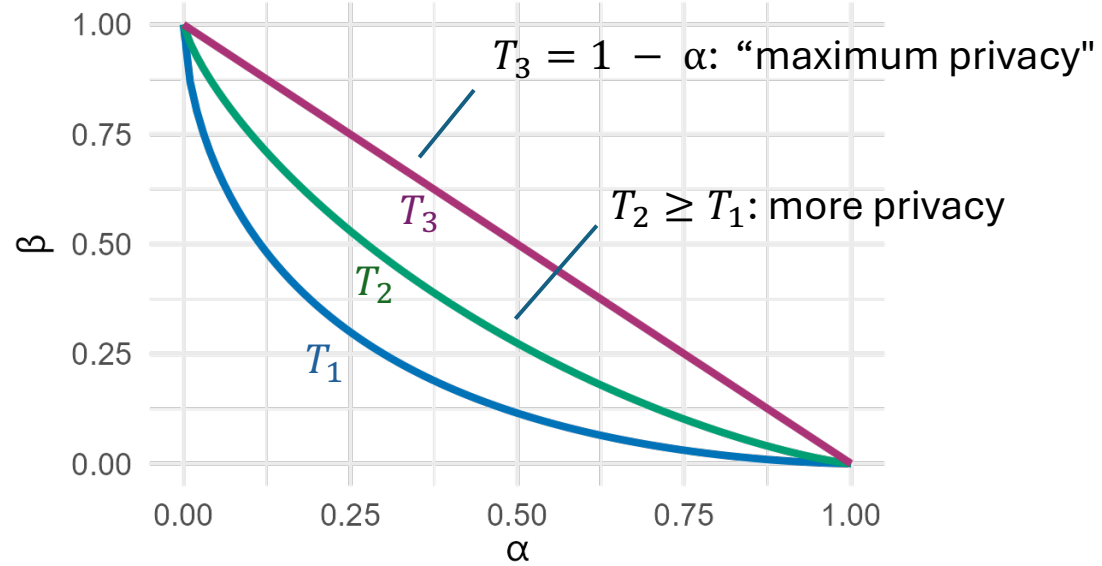


Privacy Parameters in Standard DP

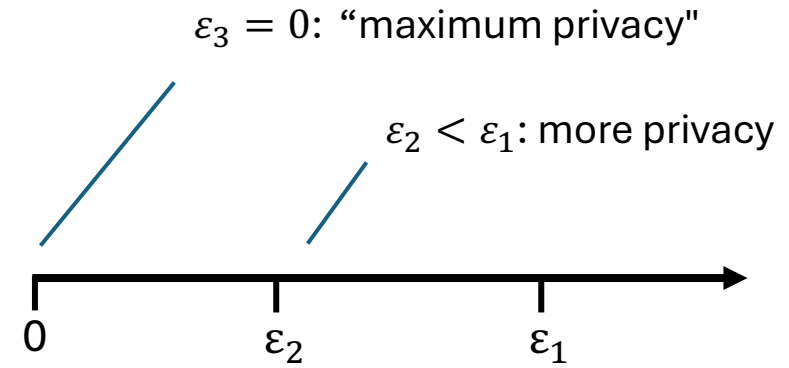


Auditing

Privacy Parameters in f -DP



Privacy Parameters in Standard DP



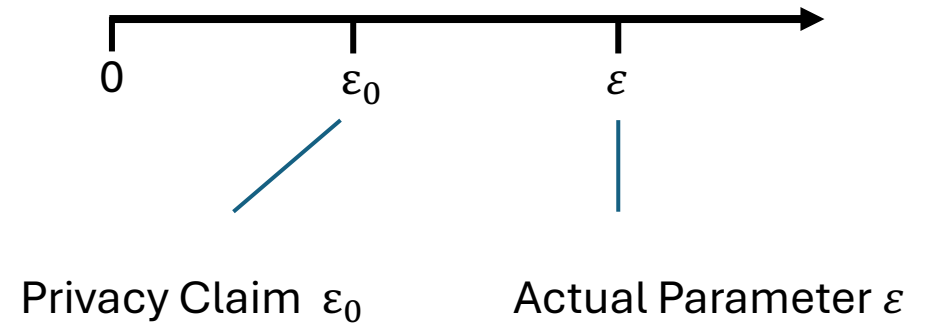
Auditing

Goal: Detect violations of privacy and expose flawed mechanisms

Auditing

Goal: Detect violations of privacy and expose flawed mechanisms

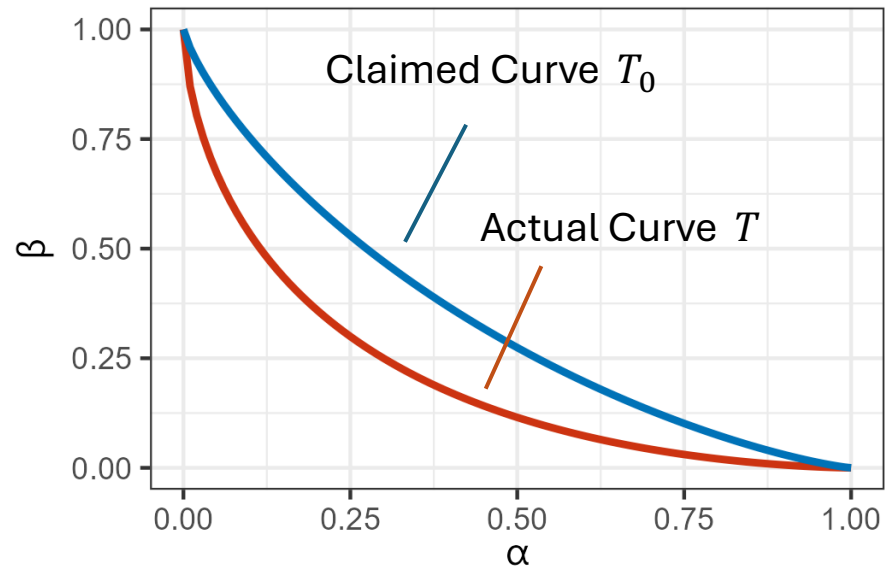
Privacy Violation in Standard DP



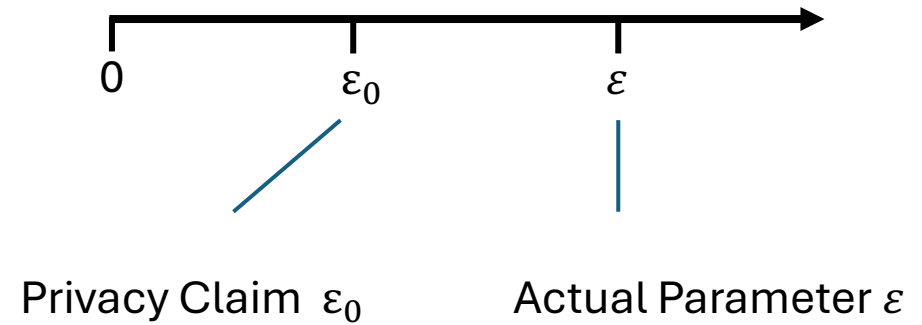
Auditing

Goal: Detect violations of privacy and expose flawed mechanisms

Privacy Violation in f -DP

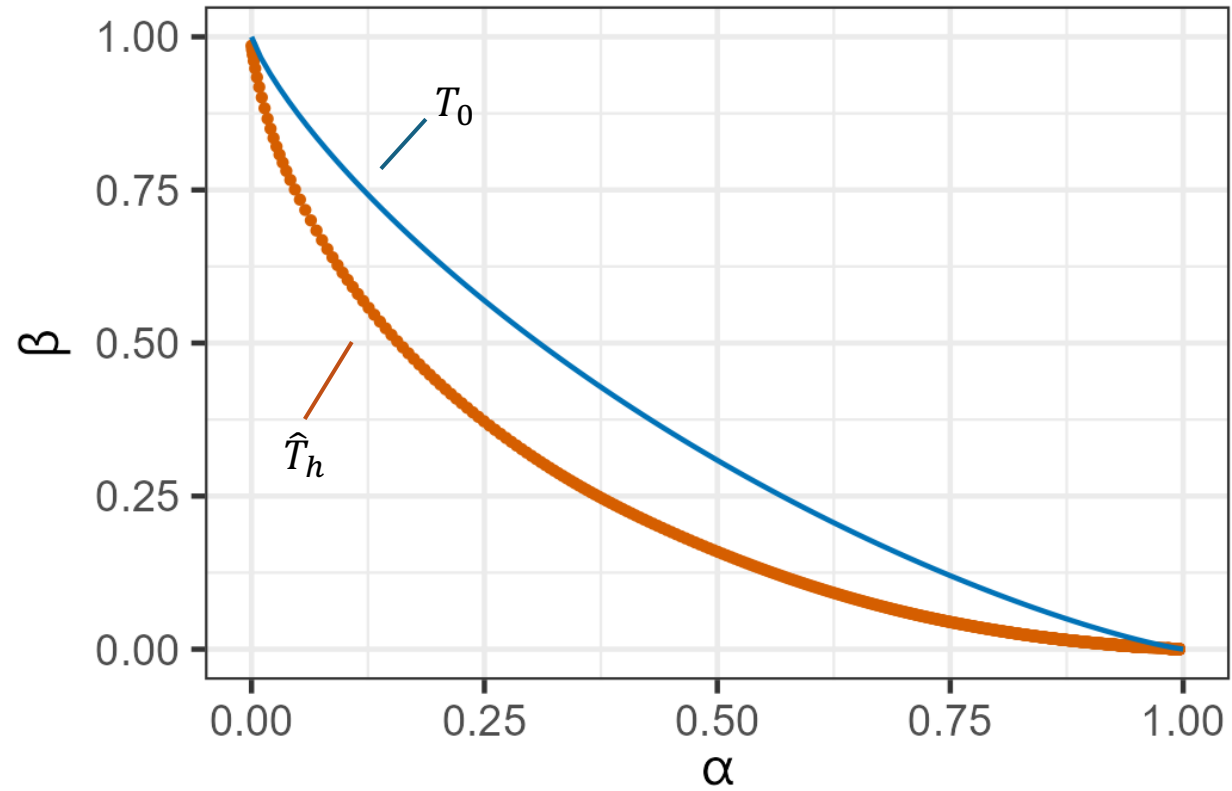


Privacy Violation in Standard DP



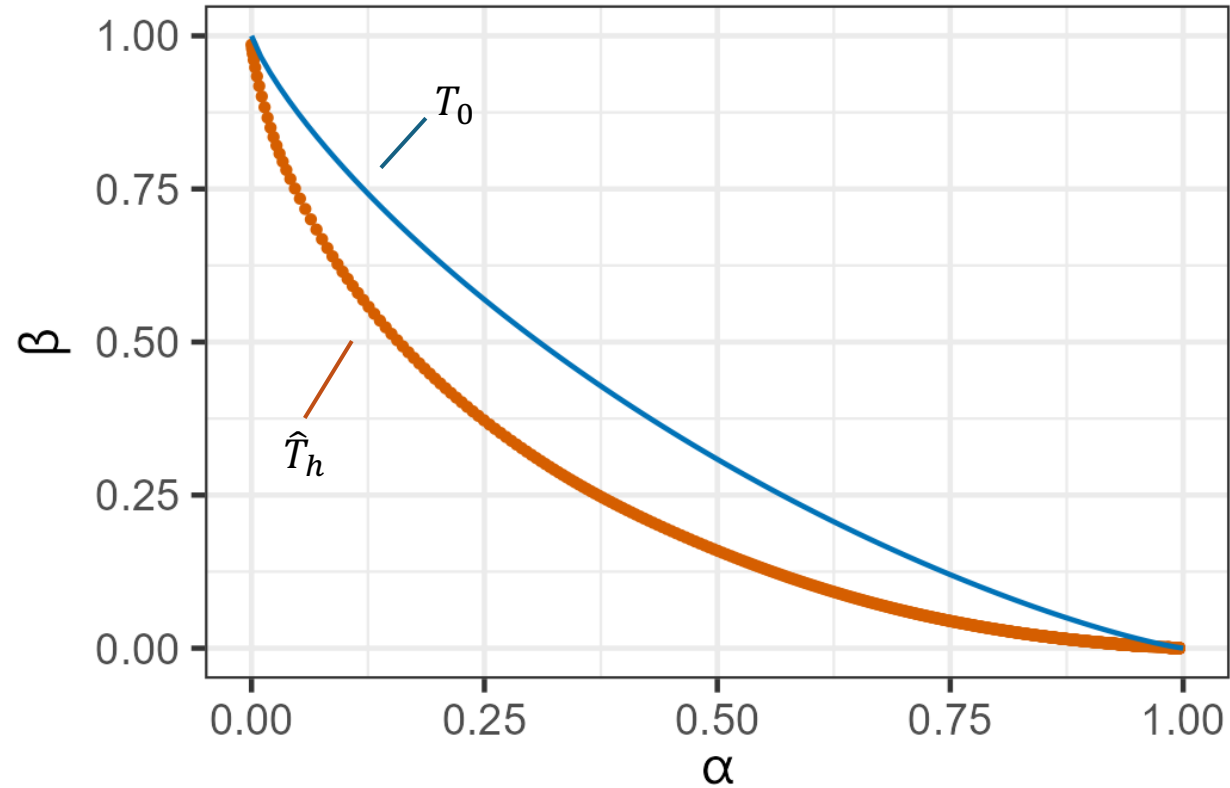
Auditing

Gaussian Mechanism



Auditing

Gaussian Mechanism



For fixed $\eta > 0$

- $\alpha(\eta)$ and $\beta(\eta)$ can be expressed as the Bayes risk of a special classification problem.
- Let $\tilde{\alpha}(\eta)$ and $\tilde{\beta}(\eta)$ be the estimates obtained via the empirical risk of the k-nearest neighbor classifier (k-NN).

Confidence Interval at η

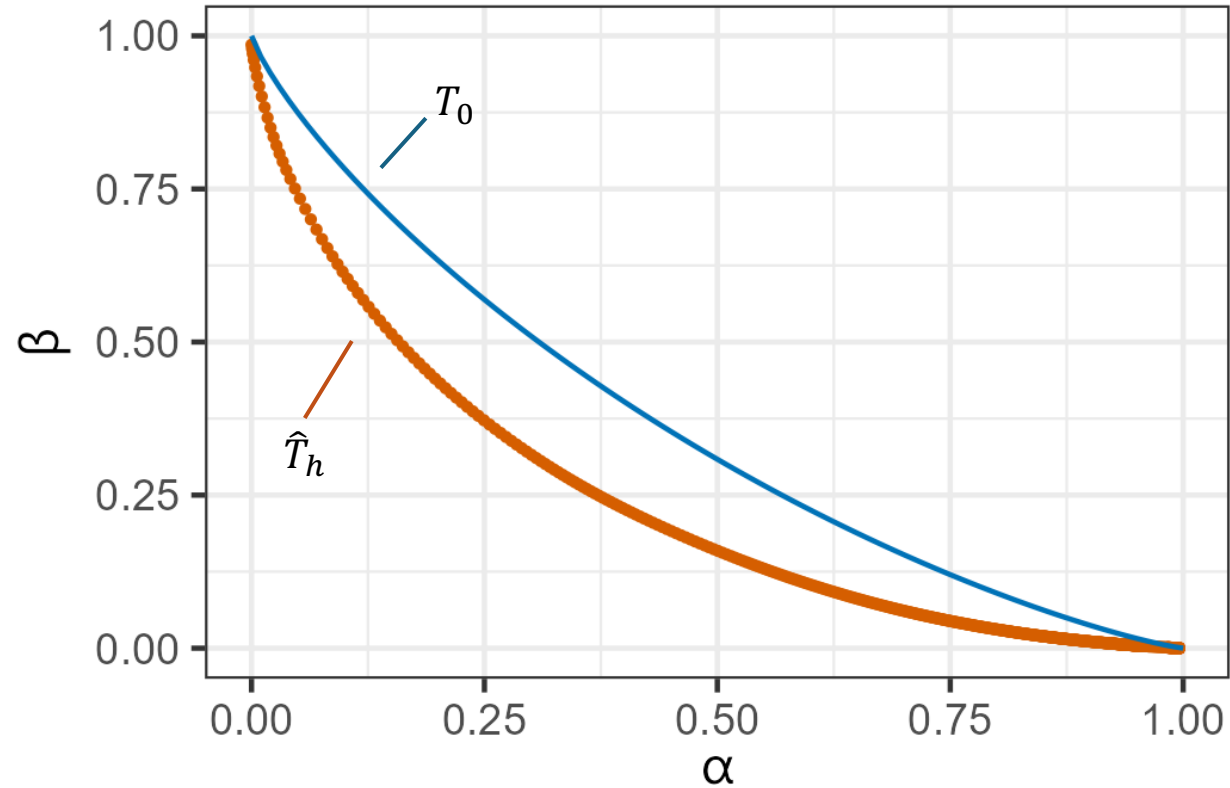
For $\gamma \in (0,1)$ and $n \geq 2$ it holds with probability greater than $1 - \gamma$ that

$$\max\{|\tilde{\alpha}(\eta) - \mathbb{E}[\tilde{\alpha}(\eta)]|, |\tilde{\beta}(\eta) - \mathbb{E}[\tilde{\beta}(\eta)]|\} \leq w(\gamma)$$

where $w(\gamma) = \sqrt{\ln(4/\gamma)} / \sqrt{2n}$.

Auditing

Gaussian Mechanism



Construct Confidence Region

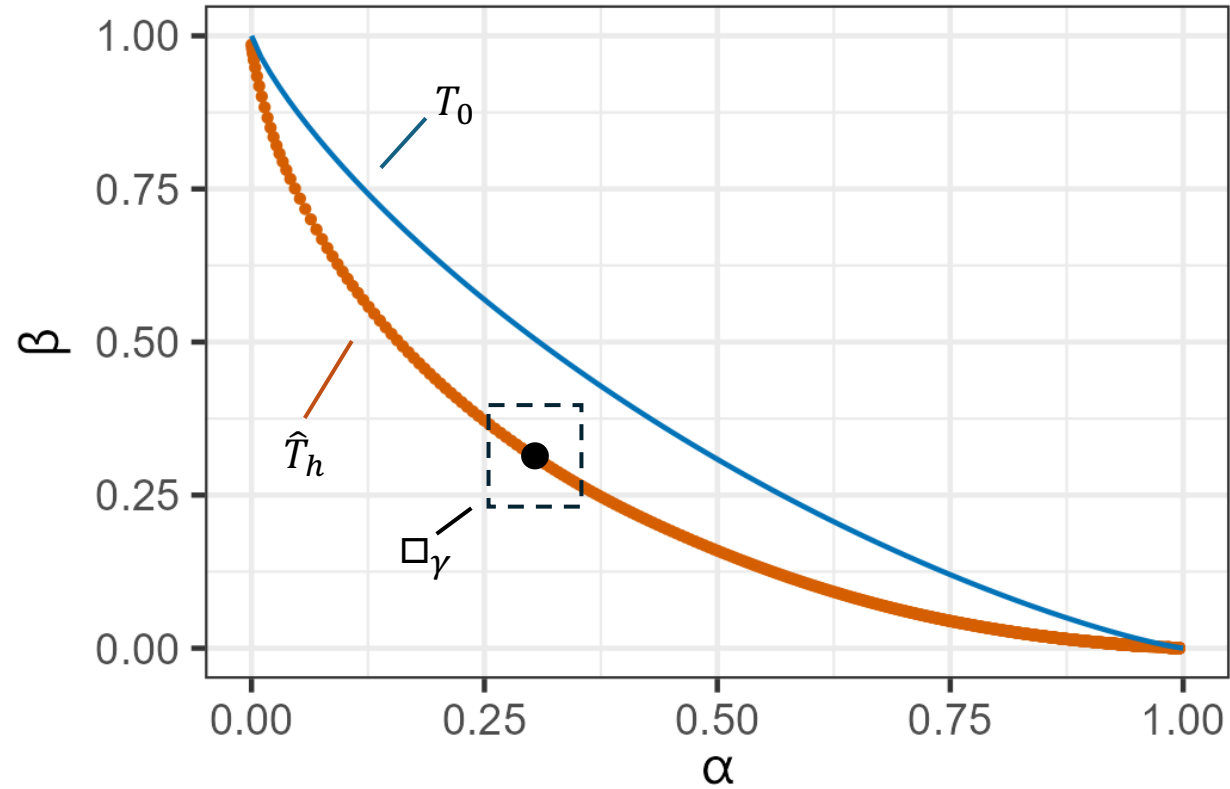
(1) Compute $\hat{\eta}^*$ where distance between $\hat{T}_h(\hat{\alpha}_h(\eta))$ and $T_0(\hat{\alpha}_h(\eta))$ is largest.

(2) Compute

$$\square_\gamma = [\tilde{\alpha}(\hat{\eta}^*) - w(\gamma), \tilde{\alpha}(\hat{\eta}^*) + w(\gamma)] \\ \times [\tilde{\beta}(\hat{\eta}^*) - w(\gamma), \tilde{\beta}(\hat{\eta}^*) + w(\gamma)].$$

Auditing

Gaussian Mechanism



Construct Confidence Region

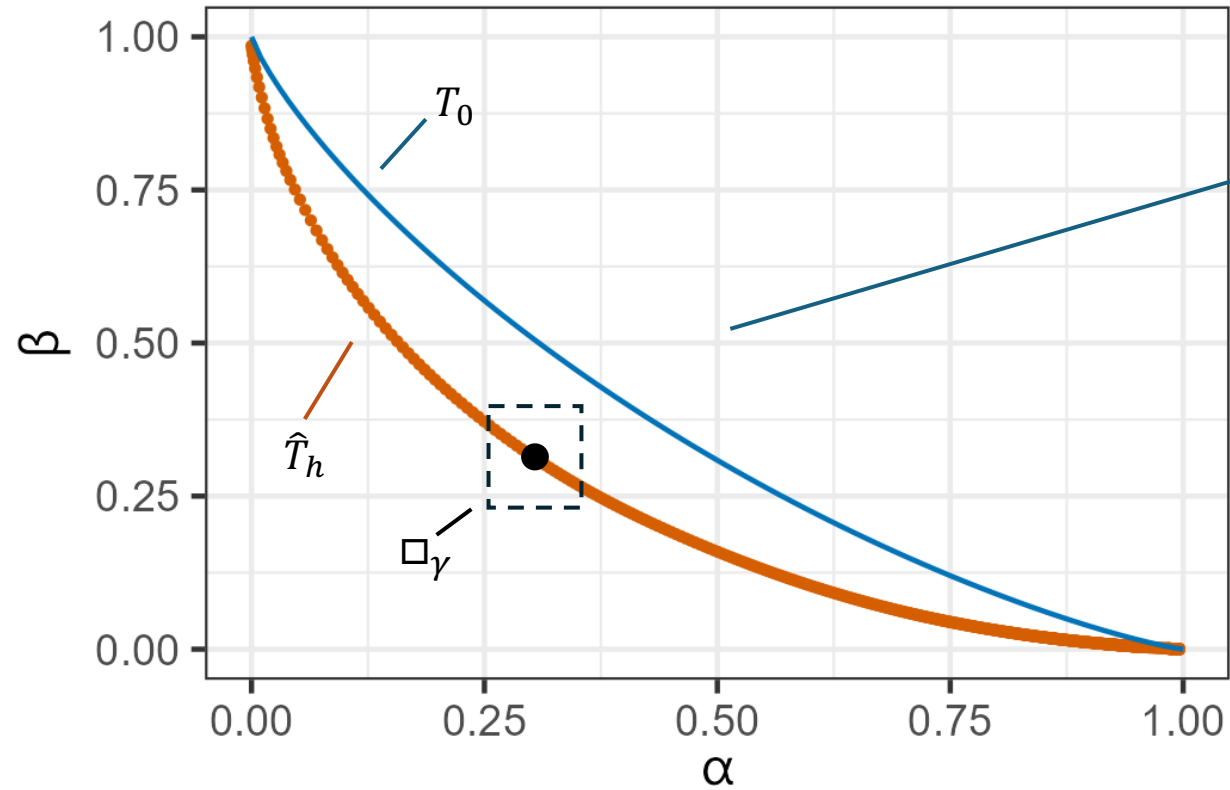
(1) Compute $\hat{\eta}^*$ where distance between $\hat{T}_h(\hat{\alpha}_h(\eta))$ and $T_0(\hat{\alpha}_h(\eta))$ is largest.

(2) Compute

$$\square_\gamma = [\tilde{\alpha}(\hat{\eta}^*) - w(\gamma), \tilde{\alpha}(\hat{\eta}^*) + w(\gamma)] \\ \times [\tilde{\beta}(\hat{\eta}^*) - w(\gamma), \tilde{\beta}(\hat{\eta}^*) + w(\gamma)].$$

Auditing

Gaussian Mechanism



Privacy Violation

$(\mathbb{E}[\tilde{\alpha}(\hat{\eta}^*)], \mathbb{E}[\tilde{\beta}(\hat{\eta}^*)])$ should lie on or above T_0 and, with high probability, in \square_γ at the same time!



Summary

- Estimation and Auditing of f -DP with Black-Box Access
- Reliable Inference and Detection of Privacy Violations
- Theoretical Guarantees

Take a look at our paper for more details on

- Algorithms
- Theory
- Experiments